

LoRA-One: One-Step Full Gradient Could Suffice for Fine-Tuning Large Language Models, Provably and Efficiently

Yuanhe Zhang^{1,2} Fanghui Liu^{2,3} Yudong Chen⁴

Abstract

This paper explores how theory can guide and enhance practical algorithms, using Low-Rank Adaptation (LoRA) (Hu et al., 2022) in large language models as a case study. We rigorously prove that, under gradient descent, **LoRA adapters align with specific singular subspaces of the one-step full fine-tuning gradient**. This result suggests that, by properly initializing the adapters using the one-step full gradient, subspace alignment can be achieved immediately—applicable to both **linear and nonlinear models**. Building on our theory, we propose a theory-driven algorithm, *LoRA-One*, where the linear convergence (as well as generalization) is built and **incorporating preconditioners theoretically helps mitigate the effects of ill-conditioning**. Besides, our theory reveals connections between *LoRA-One* and other gradient-alignment-based methods, helping to clarify misconceptions in the design of such algorithms. *LoRA-One* achieves significant empirical improvements over LoRA and its variants across benchmarks in natural language understanding, mathematical reasoning, and code generation. Code is available at: <https://github.com/YuanheZ/LoRA-One>.

1. Introduction

How to efficiently approximate or learn nonlinear models is a central question in large-scale machine learning, especially in the era of large language models (LLMs) (Brown et al., 2020; Thoppilan et al., 2022). Fine-tuning (Dodge et al., 2020) aims to make LLMs perform well on new tasks

while retain the knowledge from pre-trained models. For scalability, we expect that fine-tuning can be conducted with low computation/memory cost, i.e., parameter-efficient fine-tuning (PEFT) (Houlsby et al., 2019; Han et al., 2024).

One typical PEFT strategy is Low-Rank Adaptation (LoRA) (Hu et al., 2022), which learns an approximation of the unknown feature shift Δ by two low-rank matrices A and B with rank r , i.e. $\Delta \approx AB$ under the following initialization (denoted by index 0):

$$[A_0]_{ij} \sim \mathcal{N}(0, \alpha^2) \quad [B_0]_{ij} = 0, \quad \alpha > 0. \quad (\text{LoRA-init})$$

To improve the performance in the downstream tasks, various LoRA-based algorithms have been proposed based on, e.g., refined initialization (Li et al., 2025), learning rates (Hayou et al., 2024), efficiency (Kopiczko et al., 2024), and gradient information (Meng et al., 2024; Wang et al., 2024).

Although LoRA is conceptually simple, its optimization dynamics are inherently nonlinear and non-convex. There is few theoretical understanding of its behavior, e.g., optimization from *lazy-training* regime (Jang et al., 2024; Malladi et al., 2023; Liu et al., 2025) to *non-lazy training* regime (Kim et al., 2025) and generalization guarantees in some simplified settings (Dayi & Chen, 2024). It still remains unclear how (low-rank) gradient updates in LoRA evolve and which subspaces LoRA will converge to. More importantly, given the application-driven nature of LoRA, a rigorous theoretical understanding should not only explain its behavior but also inform practical algorithm design. **The goal of this work is to enhance LoRA’s empirical performance through theoretically grounded insights.** To this end, we address two key questions at the intersection of theory and practice:

- *Q1: How to characterize low-rank dynamics of LoRA and the associated subspace alignment in theory?*
- *Q2: How can our theoretical results contribute to algorithm design for LoRA in practice?*

1.1. Contributions

In this work, we theoretically investigate the behavior of gradient descent (GD) update of LoRA parameters (A_t, B_t) and identify the subspaces they align with. Our theory identifies the **optimal** initialization strategies in the perspective

¹Department of Statistics, University of Warwick, UK.
²Department of Computer Science, University of Warwick, UK.
³Centre for Discrete Mathematics and its Applications (DIMAP), University of Warwick, UK. ⁴Department of Computer Sciences, University of Wisconsin-Madison, USA. Correspondence to: Fanghui Liu <fanghui.liu@warwick.ac.uk>.

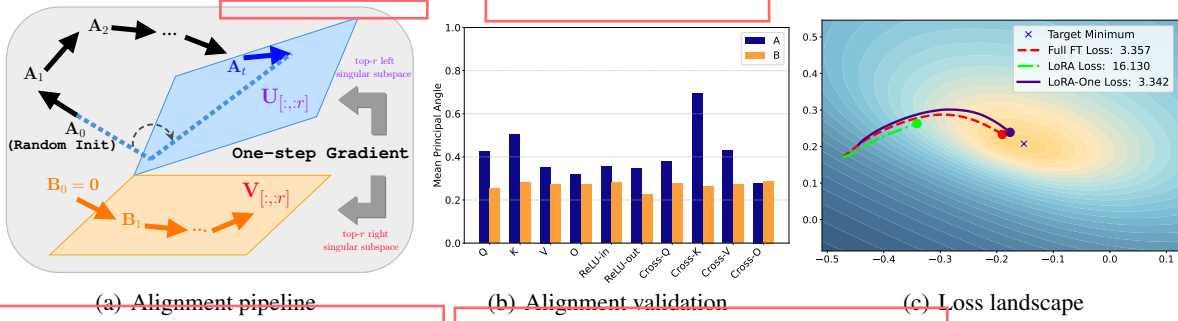


Figure 1: (a) Illustration of the alignment behavior of LoRA with the certain singular subspaces. (b) The mean principle angles within each layer class is estimated for the alignment of fine-tuning a T5 base model (Raffel et al., 2020) on MRPC using LoRA. The principal angle measures the distance between the projection matrices of the top- r (LoRA rank) singular subspace of the LoRA matrices and that of the one-step gradient. Note that the principal angles are approximately 1 at initialization and then decrease for alignment. (c) Comparison of trajectories among full fine-tuning (Full FT), LoRA, and *LoRA-One* under gradient descent. We use a two-layer neural network pretrained on odd-labeled data and then fine-tune on even-labeled MNIST data. Experimental details are presented in Appendix G.1.

of subspace alignment, and we find that it also performs well on some real-world datasets. We term this initialization as **spectral initialization**, which leverages the information of one-step full gradient, leading to the theoretical grounded algorithm, *LoRA-One*. This algorithm incorporated into several architectures achieves promising performance on natural language processing (NLP), reasoning tasks when compared to LoRA and its variants. Our contributions from theory (see Table 1 for summary) to practice are:

i) Alignment and algorithm design principles: We start by analyzing LoRA for fine-tuning a multi-output linear model. Denoting one-step gradient of full fine-tuning as G^\natural , we prove that the gradient update aligns A_t with the left top- r^* singular subspace of G^\natural while B_t always stays in a right top- r^* singular subspace w.r.t. G^\natural as shown in Fig. 1(a), where r^* is the rank of Δ . This alignment phenomenon is also empirically verified in real-world fine-tuning tasks, see Fig. 1(b). Based on the alignment results, we compute the singular value decomposition (SVD) of $G^\natural = \tilde{U}_{G^\natural} \tilde{S}_{G^\natural} \tilde{V}_{G^\natural}^\top$. The alignment can be directly achieved at the certain initialization strategy, termed as *spectral initialization*

$$\begin{aligned} A_0 &= \sqrt{\gamma} \begin{bmatrix} \tilde{U}_{G^\natural} \\ \tilde{S}_{G^\natural}^{1/2} \end{bmatrix}_{[:,1:r]} \begin{bmatrix} \tilde{S}_{G^\natural}^{1/2} \end{bmatrix}_{[1:r]}, \\ B_0 &= \sqrt{\gamma} \begin{bmatrix} \tilde{S}_{G^\natural}^{1/2} \end{bmatrix}_{[1:r]} \begin{bmatrix} \tilde{V}_{G^\natural} \end{bmatrix}_{[:,1:r]}^\top, \end{aligned} \quad (\text{Spectral-init})$$

where γ is a tuning parameter. By (Spectral-init), we theoretically ensure that $\|A_0 B_0 - \Delta\|_F$ is sufficiently small at beginning, see the theoretical results for linear models in Section 3.2 and nonlinear models in Section 4, respectively. It demonstrates the sufficiency of using one-step full gradient, which can be numerically verified on several real-world benchmarks, serving as the algorithm principle.

ii) Global convergence and generalization guarantees:

Under spectral initialization, continuing gradient descent (GD) updates for (A_t, B_t) , we further establish the linear convergence rate of $\|A_t B_t - \Delta\|_F$ for both linear and non-linear models. This linear rate, however, is sensitive to the condition number $\kappa(\Delta)$ of Δ , leading to unsatisfactory convergence performance if Δ is ill-conditioned. To address this issue, we rigorously show that adding preconditioners into the GD update eliminates the dependence on the condition number; see Section 3.2 and Section 4, respectively.

Moreover, our theory aims to clarify certain misunderstandings in prior algorithm designs. Specifically, it identifies the correct subspace for alignment and highlights potential limitations of previous LoRA variants based on gradient alignment—such as *LoRA-GA* (Wang et al., 2024); see the discussion in Section 5.

iii) Performance improvement in numerical and read-world datasets:

Guided by our theory, the spectral initialization strategy (Spectral-init) leads to our theoretically grounded algorithm, *LoRA-One*. As shown in Fig. 1(c), our numerical results demonstrate that *LoRA-One*’s trajectory is close to the full fine-tuning (Full FT) and obtain lower loss than LoRA.

We conducted experimental comparisons between *LoRA-One* and standard LoRA-based algorithms across various NLP benchmarks, including natural language understanding (NLU), mathematical reasoning, and code generation tasks. For instance, using only (Spectral-init), it takes just **one second** on some NLU tasks to achieve performance comparable to LoRA which requires tens of seconds. On the HumanEval benchmark, LLaMA 2-7B fine-tuned with *LoRA-One* achieves a score of 28.66, outperforming stan-

Table 1: Main results in the main text and appendix from subspace alignment to global convergence.

Model	Results	Algorithm	Initialization	Conclusion
Linear	Theorem 3.1	GD	(LoRA-init)	Subspace alignment of B_t
	Theorem 3.2	GD	(LoRA-init)	Subspace alignment of A_t
	Theorem 3.3	GD	(Spectral-init)	$\ A_0 B_0 - \Delta\ _F$ is small
	Theorem C.17	GD	(Spectral-init)	Linear convergence of $\ A_t B_t - \Delta\ _F$
	Theorem C.21	Precondition GD	(Spectral-init)	Linear convergence rate independent of $\kappa(\Delta)$
Nonlinear	Theorem 4.2	Precondition GD	(Spectral-init)	Linear convergence rate independent of $\kappa(\Delta)$

dard LoRA (25.85) by 2.81, while maintaining almost the same time and memory costs.

Notations For a matrix A , let $\|A\|_{op}$ denote its operator and $\|A\|_F$ its Frobenius norm. Let \odot denote the Hadamard (i.e., entrywise) matrix product. We use I_n to denote the $\mathbb{R}^{n \times n}$ -valued identity matrix. The notation U_A denotes the left singular matrix of the compact SVD of A and $U_{A,\perp}$ denotes the corresponding orthogonal complement. Similarly, V_A denotes the right singular matrix of A and $V_{A,\perp}$ denotes its orthogonal complement. Let $U_{r^*}(A)$ denote the left singular subspace spanned by the r^* largest singular values of A and $U_{r^*,\perp}(A)$ denote the left singular subspace orthogonal to $U_{r^*}(A)$. Similarly define $V_{r^*}(A)$ and $V_{r^*,\perp}(A)$ for the right singular subspace. A complete list notations can be found in Table 5 of Appendix A.

1.2. Related Work

Parameter-Efficient Fine-Tuning (PEFT): LoRA (Hu et al., 2022) and its variants have received great attention for downstream applications. The variants of LoRA focus on imbalance stepsize (Hayou et al., 2024), initialization using SVD of pre-trained weights (Meng et al., 2024), gradient approximation (Wang et al., 2024; 2025) for better performance, reducing parameters (Kopiczko et al., 2024) efficiency, preconditioned algorithm (Zhang & Pilanci, 2024) for stability.

In theory, the training dynamics and generalization ability of LoRA are rarely discovered. Based on the empirical evidence of kernel behavior of LoRA in Malladi et al. (2023), the global convergence is given by Jang et al. (2024) for LoRA with rank $\mathcal{O}(\sqrt{N})$ under the *lazy training* (Jacot et al., 2018) as well as Liu et al. (2025) on PL* condition. Beyond the lazy training regime, Kim et al. (2025) study the loss landscape of LoRA as well as its implicit bias. For generalization, Dayi & Chen (2024) derive the sample/time complexity by exploring the SGD dynamics of rank-1 LoRA, related to single-index model (Arous et al., 2021). In our work, we study the dynamics of LoRA from the perspective of subspace alignment, which has some overlap with matrix sensing as below.

Matrix Sensing under Gradient Descent: Since LoRA performs fine-tuning using a Burer-Montorio factorization, it admits similarities with matrix sensing problems, including the symmetric matrix problem with $r = r^*$ (Li et al., 2018) and $r \geq r^*$ (Stöger & Soltanolkotabi, 2021); asymmetric problem with $r \geq r^*$ (Soltanolkotabi et al., 2023; Xiong et al., 2024). Regarding initialization, small initialization (Ding et al., 2022) and spectral initialization (Ma et al., 2021) help convergence with theoretical guarantees, which is applied to LoRA under certain specific settings (Xu et al., 2025). Besides, adding preconditioner (Zhang et al., 2021; Tong et al., 2021; Xu et al., 2023; Zhang et al., 2023; Giampouras et al., 2024; Zhu et al., 2024) is beneficial to solve the problem of ill-conditioned ground truth matrix.

Technically, for the alignment part, our theory leverages some techniques from Soltanolkotabi et al. (2023). However, the symmetrization technique used in prior work cannot be applied to decouple the GD dynamics of (A_t, B_t) , posing a challenge in analyzing their individual spectral behaviors. To overcome this limitation, we develop a novel approach that enables a detailed analysis of the distinct spectral dynamics of A_t and B_t , which is one technical contribution of this work. In fact, one-step gradient information has been used in deep learning theory, demonstrating that it allows for feature learning under different stepsizes (Ba et al., 2022; Moniri et al., 2024; Cui et al., 2024; Dandi et al., 2025). Besides, for the nonlinear model part, dynamical analysis are normally based on classical gradient-based algorithm (Damian et al., 2022; Lee et al., 2024) for feature learning. Nevertheless, how such model behaves under low-rank updates under (Spectral-init) is still unclear to our knowledge.

2. Problem Settings

In this section, we introduce the problem setting of **fine-tuning pre-trained linear and nonlinear models** with the following assumptions for our theory.

2.1. Basic Assumptions

We consider both linear and nonlinear pre-trained models with multiple outputs and thus matrix parameters (instead

of vectors), which is consistent with LoRA in practice.

Assumption 2.1 (Pre-trained model). For the input $\mathbf{x} \in \mathbb{R}^d$, we denote by $\mathbf{W}^\natural \in \mathbb{R}^{d \times k}$ the **known** pre-trained parameter matrix. We assume that the pre-trained model can be linear or nonlinear with $\sigma(\cdot) = \max\{0, \cdot\}$ being the (entry-wise) ReLU activation function.

$$f_{\text{pre}}(\mathbf{x}) := \begin{cases} (\mathbf{x}^\top \mathbf{W}^\natural)^\top \in \mathbb{R}^k & \text{linear} \\ \sigma[(\mathbf{x}^\top \mathbf{W}^\natural)^\top] \in \mathbb{R}^k & \text{nonlinear} \end{cases}.$$

Note that our results can handle large dimension d and k . For fine-tuning, we assume there exists an **unknown** low-rank feature shift Δ on \mathbf{W}^\natural that we aim to estimate.

Assumption 2.2. The downstream feature matrix $\tilde{\mathbf{W}}^\natural := \mathbf{W}^\natural + \Delta$ admits an **unknown** low-rank feature shift $\Delta \in \mathbb{R}^{d \times k}$, where $\text{Rank}(\Delta) = r^* < \min\{d, k\}$.

This assumption is widely used in the literature on LoRA analysis and matrix factorization (Zhang et al., 2021; Stöger & Soltanolkotabi, 2021; Soltanolkotabi et al., 2023; Xiong et al., 2024). Next we assume the following data generation process, i.e., label-noiseless and well-behaved data.

Assumption 2.3 (Data generation process for fine-tuning). Given the unknown $\tilde{\mathbf{W}}^\natural$, the label $\tilde{\mathbf{y}}$ is generated by

$$\tilde{\mathbf{y}} := \begin{cases} (\tilde{\mathbf{x}}^\top \tilde{\mathbf{W}}^\natural)^\top \in \mathbb{R}^k, \quad \{\tilde{\mathbf{x}}_i\}_{i=1}^N \stackrel{i.i.d.}{\sim} SG, & \text{linear} \\ \sigma[(\tilde{\mathbf{x}}^\top \tilde{\mathbf{W}}^\natural)^\top], \quad \{\tilde{\mathbf{x}}_i\}_{i=1}^N \stackrel{i.i.d.}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I}_d) & \text{nonlinear} \end{cases},$$

where SG denotes the probability distribution for isotropic centered **sub-Gaussian random vectors**. We assume that we have N i.i.d training data $\{\tilde{\mathbf{x}}_i, \tilde{\mathbf{y}}_i\}_{i=1}^N$ for fine-tuning.

Note that the nonlinear model can be regarded as a special case of multi-index model (Damian et al., 2022; Abbe et al., 2022; Bietti et al., 2023) and Gaussian data is a common assumption in the analysis of single/multi-index models (Damian et al., 2022; Lee et al., 2024; Oko et al., 2024). We additionally assume that $d < N$, which coincides with practical settings of LoRA for LLaMA 2-7b (Touvron et al., 2023) on real-world datasets, e.g., MetaMathQA (Yu et al., 2024) and Code-Feedback (Zheng et al., 2024), where $d = 128 \sim 4096$ and N is on the order of 10^5 .

2.2. Full Fine-tuning and LoRA

Our goal is to efficiently recover Δ by fine-tuning on the downstream data. Let the complete SVD of $\Delta \in \mathbb{R}^{d \times k}$ be

$$\Delta = \tilde{\mathbf{U}} \tilde{\mathbf{S}}^* \tilde{\mathbf{V}}^\top := [\mathbf{U} \quad \mathbf{U}_\perp] \begin{bmatrix} \mathbf{S}^* & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{V}^\top \\ \mathbf{V}_\perp^\top \end{bmatrix}, \quad (1)$$

where $\tilde{\mathbf{U}} \in \mathbb{R}^{d \times d}$ and $\tilde{\mathbf{V}} \in \mathbb{R}^{k \times k}$ are the left and right singular matrices, and $\tilde{\mathbf{S}}^* \in \mathbb{R}^{d \times k}$ is a rank- r^* diagonal matrix with nonzero singular values $\{\lambda_i^*\}_{i=1}^{r^*}$. It admits

the compact SVD $\Delta = \mathbf{U} \mathbf{S}^* \mathbf{V}^\top$ with $\mathbf{U} \in \mathbb{R}^{d \times r^*}$, $\mathbf{V}^\top \in \mathbb{R}^{r^* \times k}$, and $\mathbf{S}^* \in \mathbb{R}^{r^* \times r^*}$. The left/right singular subspaces spanned by \mathbf{U} and \mathbf{V} play an important role in our analysis.

We write the downstream data in a compact form $\tilde{\mathbf{X}} = [\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_N]^\top \in \mathbb{R}^{N \times d}$ and the label matrix $\tilde{\mathbf{Y}} = [\tilde{\mathbf{y}}_1, \dots, \tilde{\mathbf{y}}_N]^\top \in \mathbb{R}^{N \times k}$ is generated by either linear or nonlinear target functions in Assumption 2.3. We introduce the training based on full fine-tuning and LoRA below.

Full Fine-tuning: We consider the following empirical risk minimization with a squared loss

$$L(\mathbf{W}) := \frac{1}{2N} \begin{cases} \|\tilde{\mathbf{X}} \mathbf{W} - \tilde{\mathbf{Y}}\|_{\text{F}}^2 & \text{linear,} \\ \|\sigma(\tilde{\mathbf{X}} \mathbf{W}) - \tilde{\mathbf{Y}}\|_{\text{F}}^2 & \text{nonlinear} \end{cases}, \quad (2)$$

where the parameter \mathbf{W} can be learned by gradient descent (GD) initialized at \mathbf{W}^\natural , i.e., $\mathbf{W}_0 := \mathbf{W}^\natural$.

LoRA: It updates two low-rank matrices $\mathbf{A} \in \mathbb{R}^{d \times r}$, $\mathbf{B} \in \mathbb{R}^{r \times k}$ for efficiency with the following empirical risk

$$\tilde{L}(\mathbf{A}, \mathbf{B}) := \frac{1}{2N} \begin{cases} \|\tilde{\mathbf{X}}(\mathbf{W}^\natural + \mathbf{A}\mathbf{B}) - \tilde{\mathbf{Y}}\|_{\text{F}}^2, & \text{linear,} \\ \|\sigma(\tilde{\mathbf{X}}(\mathbf{W}^\natural + \mathbf{A}\mathbf{B})) - \tilde{\mathbf{Y}}\|_{\text{F}}^2, & \text{nonlinear} \end{cases} \quad (3)$$

which can be minimized using GD with stepsize $\eta > 0$

$$\begin{aligned} \mathbf{A}_{t+1} &= \mathbf{A}_t - \eta \nabla_{\mathbf{A}} \tilde{L}(\mathbf{A}_t, \mathbf{B}_t), \\ \mathbf{B}_{t+1} &= \mathbf{B}_t - \eta \nabla_{\mathbf{B}} \tilde{L}(\mathbf{A}_t, \mathbf{B}_t). \end{aligned} \quad (4)$$

Since the true rank r^* of Δ is unknown in LoRA, our results will cover two cases: *over-ranked* ($r \geq r^*$) and *exact-ranked* ($r = r^*$).¹ Our results allow for large d, k while $r, r^* = \Theta(1)$, which coincides with common practice.

Optimization and Generalization: We are interested in the error $\|\mathbf{A}_t \mathbf{B}_t - \Delta\|_{\text{F}}^2$ under the LoRA training dynamics. Bounds on this error also imply generalization performance, because the generalization error for a new data $(\tilde{\mathbf{x}}, \tilde{\mathbf{y}})$ satisfies $\mathbb{E}_{\tilde{\mathbf{x}}} \|\tilde{\mathbf{y}} - \sigma(\mathbf{W}^\natural + \mathbf{A}_t \mathbf{B}_t)^\top \tilde{\mathbf{x}}\|_2^2 \leq \|\mathbf{A}_t \mathbf{B}_t - \Delta\|_{\text{F}}^2$ in the nonlinear setting, with equality in the linear setting.

3. Analysis of LoRA under Linear Model

In this section, we establish the alignment between LoRA and one gradient of full fine-tuning. This result guides us to design new strategies for speeding up practical LoRA-based algorithms, which achieve this alignment at initialization.

We formally define the negative gradient of full fine-tuning

¹In the matrix sensing/completion literature, they are often called *over-* and *exact-parameterized*, respectively.

in Eq. (2) for the linear setting after the first step as

$$\mathbf{G}^{\natural} := -\nabla_{\mathbf{W}} L(\mathbf{W}^{\natural}) = \frac{1}{N} \widetilde{\mathbf{X}}^{\top} (\widetilde{\mathbf{Y}} - \widetilde{\mathbf{X}} \mathbf{W}^{\natural}). \quad (5)$$

Note that $\widetilde{\mathbf{X}}^{\top} \widetilde{\mathbf{X}}$ is a non-singular square matrix (Zeng & Lee, 2024, Lemma 6). Since left multiplication by a non-singular square matrix does not change the rank (Horn & Johnson, 2012, 0.4.6 (b)), we have $\text{Rank}(\mathbf{G}^{\natural}) = \text{Rank}(\Delta) = r^*$. Then, we denote the singular values of \mathbf{G}^{\natural} by $\{\lambda_i(\mathbf{G}^{\natural})\}_{i=1}^{r^*}$ in non-increasing order.

3.1. Alignment under LoRA Initialization

We first present the results for the alignment of \mathbf{B}_t by recalling the notations $\mathbf{V}_{r^*}(\cdot)$ and $\mathbf{V}_{r^*,\perp}(\cdot)$.

Theorem 3.1 (Alignment between \mathbf{G}^{\natural} and \mathbf{B}_t). *Under assumptions in Section 2.1 for the linear setting, consider the LoRA updates (4) with (LoRA-init). We have*

$$\|\mathbf{V}_{r^*,\perp}^{\top}(\mathbf{G}^{\natural}) \mathbf{V}_{r^*}(\mathbf{B}_t)\|_{op} = 0, \quad \forall t \in \mathbb{N}_+.$$

One can see that, due to the zero initialization of \mathbf{B}_0 in (LoRA-init), after the first GD step, it holds that $\mathbf{B}_1 = \eta \mathbf{A}_0^{\top} \mathbf{G}^{\natural}$, which has rank $\leq r^*$ and lies in the right top- r^* singular subspace of \mathbf{G}^{\natural} . The subsequent GD dynamics of \mathbf{B}_t is always restricted to this invariant subspace.

Next we build the alignment for \mathbf{A}_t with the notations $\mathbf{U}_{r^*}(\cdot)$, $\mathbf{U}_{r^*,\perp}(\cdot)$ and κ^{\natural} as the condition number of \mathbf{G}^{\natural} .

Theorem 3.2 (Alignment between \mathbf{G}^{\natural} and \mathbf{A}_t . Simplified version of Theorem C.9). *For the $r \geq 2r^*$ case, under assumptions in Section 2.1 for the linear setting, we consider the LoRA updates (4) with $[\mathbf{A}_0]_{ij} \sim \mathcal{N}(0, \alpha^2)$ in (LoRA-init). Then for any constant $\theta \in (0, 1)$, by taking $\alpha = \mathcal{O}\left(\theta^{\frac{3}{2}\kappa^{\natural}} d^{-\frac{3}{4}\kappa^{\natural} - \frac{1}{2}} \|\mathbf{G}^{\natural}\|_{op}^{\frac{1}{2}}\right)$, and running gradient descent for t^* steps with*

$$t^* \lesssim \frac{\ln\left(\frac{\sqrt{d}}{\theta}\right)}{\ln(1 + \eta \lambda_{r^*}(\mathbf{G}^{\natural}))}, \quad (6)$$

we achieve the following the alignment on the left singular subspace between \mathbf{G}^{\natural} and \mathbf{A}_{t^*} as below

$$\|\mathbf{U}_{r^*,\perp}^{\top}(\mathbf{G}^{\natural}) \mathbf{U}_{r^*}(\mathbf{A}_{t^*})\|_{op} \lesssim \theta, \quad (7)$$

with probability at least $1 - C_1 \exp(-d) - C_2 \exp(-r) - C_3 \exp(-N)$ for some constants C_1, C_2, C_3 .

Remark: The result under the $r^* \leq r < 2r^*$ case is more complex and we defer this result to Theorem C.9. The choice of α in Theorem 3.2 shows that after $t^* = \Theta\left(\frac{\ln d}{\lambda_{r^*}(\mathbf{G}^{\natural})}\right)$ in Eq. (6), the alignment can be achieved. Our results can cover the standard He-initialization (He et al., 2015) if $\|\mathbf{G}^{\natural}\|_{op} \geq \Omega(d^{\frac{3}{4}\kappa^{\natural}})$.

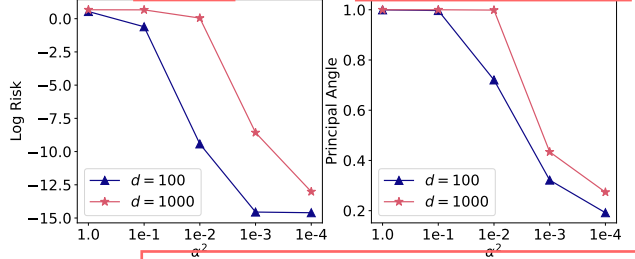


Figure 2: Under (LoRA-init), the risk and the alignment to full one-step GD of LoRA with different α^2 and d , trained via GD on task (3). *Left*: the log risk under different initialization variance α^2 . The risk is defined as $\frac{1}{2} \|\mathbf{A}_t \mathbf{B}_t - \Delta\|_F^2$. *Right*: the best principal angle between the top- r singular subspace of \mathbf{G}^{\natural} and \mathbf{A}_t during training. Smaller is closer. The principal angle is defined as $\min_t \|\mathbf{U}_{r^*,\perp}^{\top}(\mathbf{G}^{\natural}) \mathbf{U}_{r^*}(\mathbf{A}_t)\|_{op}$. More experimental details can be found in Appendix G.1.

Requirement on $\|\mathbf{G}^{\natural}\|_{op}$ can be relaxed under smaller initialization, illustrated by Fig. 2.

The above two theorems characterize the alignment between \mathbf{G}^{\natural} and $(\mathbf{A}_t, \mathbf{B}_t)$. Fig. 2 empirically validates Theorem 3.2 in two folds: *i*) Smaller initialization (α^2 in the x-axis) encourages better alignment (evaluated by the principal angle), and then better generalization performance of fine-tuning (evaluated by the risk). But in practice smaller initialization would increase the training time for convergence, as a double-edge sword. *ii*) increasing d leads to longer alignment time, illustrated by Eq. (6), and worse alignment performance, illustrated by the formulation of α . Besides, we also verify this alignment in read-world applications by fine-tuning a T5 base model (Raffel et al., 2020) on MRPC using LoRA, as shown in Fig. 1(b). The mean principle angles within each layer class is computed in a similar way of Fig. 2. Our empirical results demonstrate that the principal angles decrease from around 1 because of Gaussian initialization to the value around 0.2 ~ 0.4.

Proof of Sketch: Here we give a proof of sketch of Theorem 3.2. The dynamics (4) can be written as

$$\underbrace{\begin{bmatrix} \mathbf{A}_{t+1} \\ \mathbf{B}_{t+1}^{\top} \end{bmatrix}}_{=: \mathbf{Z}_{t+1}} = \underbrace{\begin{bmatrix} \mathbf{I}_d & \eta \mathbf{G}^{\natural} \\ \eta \mathbf{G}^{\natural\top} & \mathbf{I}_k \end{bmatrix}}_{=: \mathbf{H}} \underbrace{\begin{bmatrix} \mathbf{A}_t \\ \mathbf{B}_t^{\top} \end{bmatrix}}_{=: \mathbf{Z}_t} + \text{nonlinear term},$$

where \mathbf{H} is a time-independent matrix corresponding to the linear part of the dynamic $\mathbf{Z}_t^{1in} := \mathbf{H}^t \mathbf{Z}_0$. By Schur decomposition of \mathbf{H} (see Lemma C.1), we can obtain the precise spectral dynamics of \mathbf{Z}_t^{1in} and derive the alignment between \mathbf{Z}_t^{1in} and \mathbf{G}^{\natural} , see Lemma C.5 for details. We prove that the nonlinear term is well controlled, i.e., $\|\mathbf{Z}_t - \mathbf{Z}_t^{1in}\|_{op} \leq \|\mathbf{A}_0\|_{op}, \forall t \leq t^*$, see Lemma C.6. Then

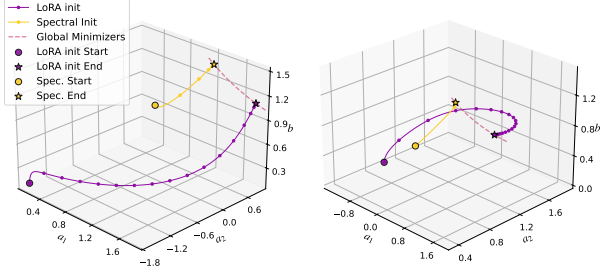


Figure 3: Comparison of the GD trajectories under (Spectral-init) and (LoRA-init) with two different starting points. See more experimental details in Appendix G.1.

the alignment between Z_t^{lin} and G^{h} can be successfully transferred to that of Z_t , see Theorem C.9 for details.

We remark that previous work on matrix sensing (Stöger & Soltanolkotabi, 2021; Soltanolkotabi et al., 2023) via a symmetrization technique cannot be directly applied to our setting. Such symmetrization technique prevents the alignment results decoupling into two factorized matrices. We extend their technique to decouple the alignment for A_t and B_t individually via Schur decomposition of H .

3.2. Spectral Initialization and Global Convergence

Theorem 3.2 has demonstrated the alignment on the rank- r^* singular space of G^{h} and (A_t, B_t) . In other words, if we take the SVD of G^{h} and choose the certain singular subspace for initialization in (Spectral-init), we can directly achieve the alignment at this initialization and recover Δ to some extent, which is the main target of this work.

By the following standard concentration result for (sub)-Gaussian data: with probability at least $1 - 2C \exp(-\epsilon^2 N)$ for some constants $C > 0$, we have

$$\left\| \widehat{\Sigma} - I_d \right\|_{\text{op}} \leq \epsilon := \min \left\{ \frac{1}{2\kappa}, \frac{c}{\kappa^3} \right\} \leq \frac{1}{2}. \quad (8)$$

Recall κ is the condition number of Δ and $\lambda_{r^*}^*$ is the r^* -th singular value of Δ , we have the following result at the spectral initialization.

Theorem 3.3. *[One-step gradient can suffice] Under assumptions in Section 2.1 for the linear setting via (Spectral-init), taking ϵ in Eq. (8), then with probability at least $1 - 2C \exp(-\epsilon^2 N)$ for constant $C > 0$, we have*

$$\|A_0 B_0 - \Delta\|_{\text{op}} \leq \epsilon \|\Delta\|_{\text{op}} \leq \frac{\lambda_{r^*}^*}{2}.$$

Theorem 3.3 demonstrates that, after one-step full gradient, i.e., using spectral initialization (Spectral-init), $A_0 B_0$ is able to recover Δ with small error. This is still true for

nonlinear models, see Lemma D.5 for details. Fig. 3 numerically validates that the starting points initialized by (Spectral-init) are consistently closer to the set of global minimizers, whereas those initialized by (LoRA-init) tend to be farther away across different random seeds.

Moreover, running gradient descent from points initialized by (Spectral-init) requires significantly fewer steps to reach a global minimizer, demonstrating the advantages of (Spectral-init). Due to page limit, we present the global convergence of linear models is deferred to Appendix C.2 and Appendix C.3, respectively. In Appendix C.2, we derive the linear convergence rate $\|A_t B_t - \Delta\|_{\text{F}}$. Since the convergence will be slow if the downstream feature shift Δ is ill-conditioned, i.e., κ is large. This motivates us to add preconditioners, then the convergence rate will be independent of κ accordingly, see Appendix C.3 for details.

4. Analysis of LoRA under Nonlinear Models

Now we focus on the nonlinear setting described in Section 2, where we consider the exact-rank case $r = r^*$ for delivery. We will demonstrate that $\|A_0 B_0 - \Delta\|_{\text{F}}$ is still small under the spectral initialization. Besides, the linear convergence rate of $\|A_t B_t - \Delta\|_{\text{F}}$ can still hold.

As an example, we demonstrate the equipment of precondition GD on (A_t, B_t) for global convergence

$$\begin{aligned} A_{t+1} &= A_t - \eta \nabla_A \tilde{L}(A_t, B_t) (B_t B_t^\top)^{-1}, \\ B_{t+1} &= B_t - \eta (A_t^\top A_t)^{-1} \nabla_B \tilde{L}(A_t, B_t). \end{aligned} \quad (9)$$

Notice that here we use standard matrix inversion since we can prove that A_t and B_t stay non-singular across all $t \geq 0$. By denoting $W_t := W^{\text{h}} + A_t B_t$, we have the gradient

$$\nabla_A \tilde{L}(A_t, B_t) = -J_{W_t} B_t^\top, \nabla_B \tilde{L}(A_t, B_t) = -A_t^\top J_{W_t},$$

where we denote

$$J_{W_t} := \frac{1}{N} \tilde{X}^\top \left[\sigma(\tilde{X} \tilde{W}^{\text{h}}) - \sigma(\tilde{X} W_t) \right] \odot \sigma'(\tilde{X} W_t).$$

To deliver the proof, apart from the above-mentioned assumptions in Section 2.1 for the the nonlinear setting, we also need the following assumption.

Assumption 4.1. We assume that i) $\frac{\|\tilde{W}^{\text{h}}\|_{\text{op}}}{\|\tilde{w}_m^{\text{h}}\|_2} = \mathcal{O}(1)$; ii) $\frac{\max\{\lambda_{r^*}^*, \|\Delta_m\|_{\text{op}}\}}{\|\tilde{w}_m^{\text{h}}\|_2} = \mathcal{O}\left(\frac{1}{\kappa r^*}\right)$ for $m \in [k]$.

Remark: The condition i) ensures the balance between different neurons within one layer for the downstream teacher model and the task diversity. The condition ii) ensures the signal of downstream feature shift is smaller than the pre-trained ones approximately in the order $(\kappa r^*)^{-1}$ since the

signal of adapted weight is generally weaker than the pre-trained weight. Two conditions can be empirically observed in Appendix G.5.

Here we can show that, for the nonlinear model, LoRA training can achieve global linear convergence under (Spectral-init) via preconditioned GD in Eq. (9).

Theorem 4.2 (Simplified version of Theorem D.10). *Under assumptions in Section 2.1 for the nonlinear setting and 4.1, with training conducted by Eq. (9) and initialization via (Spectral-init) with setting $\gamma = 2$, we take $\epsilon = \mathcal{O}\left(\frac{1}{r^* \kappa \sqrt{d}}\right)$ and $\rho \leq \frac{1}{20}$. Then choosing $\eta \in (c_\eta, 1)$ for a small constant $c_\eta > 0$, with probability at least $1 - 2Cdk \exp(-\epsilon^2 N)$ for a universal constant $C > 0$, we have*

$$\|A_t B_t - \Delta\|_F \leq \left(1 - \frac{\eta}{4}\right)^t \rho \lambda_{r^*}^*, \forall t \geq 0. \quad (10)$$

Remark: We make three remarks here:

- i) This theorem is based on $\|A_0 B_0 - \Delta\|_F \leq \rho \lambda_{r^*}^*$ at initialization, see Lemma D.5 for details, which demonstrates that one-step full gradient can be sufficient.
- ii) The convergence rate is independent of condition number κ of downstream feature shift Δ , demonstrating the benefits of adding preconditioners.

Proof of Sketch The complete proof can be found in Appendix D.2. We first compute the expectation of J_{W_t} (see Lemma D.2) and decompose J_{W_t} into $\frac{1}{2}(A_t B_t - \Delta) + \Xi_t$, where Ξ_t is defined in Lemma D.6. The first term is the signal term which can dominate the preconditioned GD dynamics. The second term $\Xi_t := T1 + T2$ consists of two parts (details see Lemma D.7): the first part $T1$ is the residual term from $\mathbb{E}_{\tilde{x}}[J_{W_t}]$ which vanishes due to pre-training signal dominance. For the second term $T2$, it comes from the concentration error of J_{W_t} , which can also controlled by large sample size N .

To handle $\|A_t B_t - \Delta\|_F$, we explore its recursion relationship in Lemma D.6. The key part is to control $\|(I_d - U_{A_t} U_{A_t}^\top) \Delta (I_k - V_{B_t} V_{B_t}^\top)\|_F$ (Lemma D.9) and higher order term (Lemma D.8).

5. Algorithm and Discussions

In this section, we present the *LoRA-One* algorithm and justify the optimality of our initialization over previous gradient alignment based algorithms for fine-tuning.

We present the implementations in Algorithm 1, which is driven by (Spectral-init) (shown in line 3-6). It coincides with the spirit of gradient alignment work, e.g., *LoRA-GA* (Wang et al., 2024), *LoRA-pro* (Wang et al., 2025), but the mechanisms for gradient alignment differ significantly, as suggested by our theory. First, *LoRA-GA* proposes the fol-

Algorithm 1 LoRA-One for one specific layer

Input: Pre-trained weight W^h , batched data $\{\mathcal{D}_m\}_{m=1}^T$, sampled batch data \mathcal{B} , LoRA rank r , LoRA alpha α , loss function L , scaling parameter s

Initialize:

- 1: Compute $\nabla_W L(W^h)$ given \mathcal{B}
- 2: $U, S, V \leftarrow \text{SVD}(-\nabla_W L(W^h))$
- 3: $S \leftarrow S/S_{[0,0]}$ and $\gamma \leftarrow 1/s$
- 4: $A_0 \leftarrow \sqrt{\gamma} \cdot U_{[:,1:r]} S_{[r,r]}^{1/2}$
- 5: $B_0 \leftarrow \sqrt{\gamma} \cdot S_{[r,r]}^{1/2} V_{[:,1:r]}^\top$
- 6: Clear $\nabla_W L(W^h)$

Train:

- 7: **for** $t = 0, \dots, T-1$ **do**
- 8: Compute gradients given \mathcal{D}_{t+1} :
 $G_{t+1}^A \leftarrow \nabla_A \tilde{L}(A_t, B_t), G_{t+1}^B \leftarrow \nabla_B \tilde{L}(A_t, B_t)$
- 9: Update $A_{t+1}, B_{t+1} \leftarrow \text{AdamW}(G_{t+1}^A, G_{t+1}^B)$
- 10: **end for**

Return: $W^h + \frac{\alpha}{\sqrt{r}} A_T B_T$

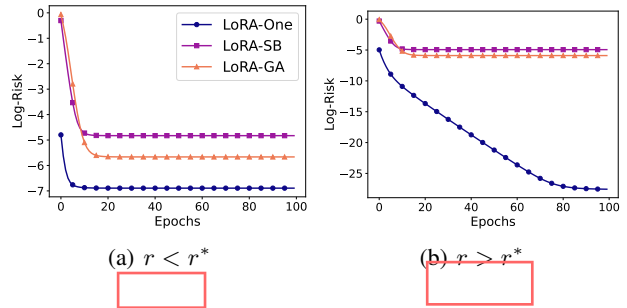


Figure 4: The log-risk curve under *LoRA-One*, *LoRA-SB*, and *LoRA-GA*, trained via GD on fine-tuning task (3) under: 1) under-ranked case $r < r^*$, 2) over-ranked case $r > r^*$.

lowing initialization strategy (omit the scaling parameters)

$$A_0 \leftarrow -[\tilde{U}_{G^h}]_{[:,1:r]}, B_0 \leftarrow [\tilde{V}_{G^h}]_{[:,r+1:2r]}^\top,$$

which aims to provide the best $2r$ approximation of G^h . However, our theory indicates that B_t will align to the right-side rank- r^* singular subspace of G^h under random initialization. However, *LoRA-GA* chooses the $(r+1)$ -th to $2r$ -th singular values for B_0 , causing the iterates B_t to lie outside the desired subspace. As a result, the optimization may remain trapped in an undesirable subspace and fail to converge to an optimal solution, which can numerically verified by Fig. 4. Moreover, this approach subtracts the gradient for non-zero initialization and thus yields a biased estimate of Δ , scaling with the model size; see further discussion in Appendix F.

Secondly, there is one concurrent work (Ponkshe et al., 2024), *LoRA-SB*, which uses the same singular subspace

~~LoRA-SB: G SVD, $r \times r$ matrix R from the SVD of G^\dagger . Their intuition is to project the fine-tuning updates onto the singular subspace of first gradient step G^\dagger . However, the singular subspace of G^\dagger normally still have distances with the ground truth, which will make their method hard to escape/rotate this subspace with limited degrees of freedom.~~

Our toy experiments based on (3) in Fig. 4 show that both *LoRA-GA* and *LoRA-SB* fail to find global minimizers even in a simple linear setting, whereas *LoRA-One* demonstrates significantly better generalization. More experimental details are presented in Appendix G.1.

6. Experiments

In this section, we conduct experiments to compare *LoRA-One* with typical LoRA based algorithms across multiple NLP benchmarks. In Section 6.1, we evaluate the ability of one-step gradient on real-world fine-tuning tasks to justify our theory on natural language understanding, i.e. Theorem 3.3 and Lemma D.5. In Section 6.2, we evaluate on mathematical reasoning, general knowledge, and code generation tasks, with more data and epochs for further evaluating math reasoning ability in Section 6.3. Furthermore, we compare the time and memory cost across different methods to illustrate our efficiency.

6.1. One-Step Full Gradient Could Suffice in Natural Language Understanding

We fine-tune T5 base model (Raffel et al., 2020) on a subset from GLUE (Wang et al., 2019) - MNLI, SST2, CoLA, QNLI, and MRPC. We evaluate the test performance by accuracy (%). We compare LoRA (Hu et al., 2022), *LoRA+* (Hayou et al., 2024), *P-LoRA* (Zhang & Pilanci, 2024), *PiSSA* (Meng et al., 2024), *LoRA-GA* (Wang et al., 2024), *LoRA-Pro* (Wang et al., 2025), and *LoRA-One* with rank 8. The hyperparameters are optimized for each method. More experimental details are presented in Appendix G.4.

Before experimental comparison, we first access the capacity of the one-step full gradient with its low-rank components on these real-world fine-tuning tasks. We approximate the one-step full-batch update G^\dagger from full fine-tuning by a large sampled batch (2048) as G_B^\dagger and a best rank- r approximation of G^\dagger for $r=8$ using a smaller sampled batch (8) as $\mathcal{P}_r(G_B^\dagger)$. We optimize the learning rate, i.e. η^* , and update the base model by $W^\dagger - \eta^* G_B^\dagger$ and $W^\dagger - \eta^* \mathcal{P}_r(G_B^\dagger)$ for SST2, CoLA, QNLI, and MRPC, respectively.

The top rows of Table 2 shows that the test performance can be significantly improved over the pre-trained model by the one-step full gradient step with proper selection of stepsize. This improvement is still promising (even better) after taking the best rank- r approximation with smaller

sampled batch, which is equivalent to (Spectral-init). We remark that low-rank update with small batch for CoLA and MRPC only costs *less than one second* but already matches the performance of LoRA which needs tens of seconds. Accordingly, **one-step full gradient can suffice for fine-tuning on small-scale datasets**, e.g., CoLA, MRPC.

Besides, Table 2 also shows that *LoRA-One* outperforms other LoRA-based methods on three tasks out of five and achieves the best in average. Significant gains appear on the smaller benchmarks, i.e. CoLA and MRPC. On large datasets such as MNLI and QNLI, *LoRA-Pro* performs better but is with more cost. This is because, *LoRA-pro* (Wang et al., 2025) approximates gradient from full fine-tuning at every training step while *LoRA-One* only conducts at the first step. *LoRA-pro* adds $10dr^2 + 6kr^2 + 155r^3/6$ more FLOPs than *LoRA-One* per pair of matrices for time cost. For memory cost, *LoRA-pro* on MetaMathQA100k with rank 8 costs 43.87 GB while our method only costs 21.7 GB for memory management.

6.2. Natural Language Generation

We fine-tune LLaMA 2-7B (Touvron et al., 2023) on: 1) 100K samples from MetaMathQA (Yu et al., 2024) and evaluate the accuracy based on two types of prompting: (a) direct prompting, (b) 8-shot Chain-of-Thought² (CoT) (Wei et al., 2022) prompting on GSM8K (Cobbe et al., 2021); 2) Alpaca (Taori et al., 2023) and evaluate on the MMLU (Hendrycks et al., 2021) benchmarks using direct prompting; 3) 100K samples from Code-Feedback (Zheng et al., 2024) and evaluate the PASS@1 on HumanEval (Chen et al., 2021a). We compare LoRA, *LoRA-GA*, and *LoRA-One* with rank 8. The learning rate and batch size are optimized for each method. More experimental details are presented in Appendix G.2.

Table 3 shows that *LoRA-One* consistently outperforms both vanilla LoRA and *LoRA-GA* across different tasks and prompting methods. *LoRA-One* achieves 60.44% accuracy under direct prompting—about 1.18 points higher than LoRA—and 55.88% in the few-shot CoT setting, a gain of roughly 2.52 points, indicating it not only strengthens the model’s core problem-solving abilities but also its capacity for coherent, multi-step reasoning. On MMLU, *LoRA-One* also shows superior generalization on the MMLU benchmark (47.24% vs. 45.73% for LoRA), indicating improved knowledge retention across diverse domains. Finally, *LoRA-One* excels in code generation, it achieves a PASS@1 score of 28.66%, nearly 3 points higher than LoRA’s 25.85% and also improving upon *LoRA-GA*, implying better adaptation to structured **code synthesis tasks**. Moreover, *LoRA-One* exhibits noticeably lower run-to-run variability compared to

²<https://github.com/EleutherAI/lm-evaluation-harness>

Table 2: Accuracy comparison on GLUE subset across typical LoRA based algorithms, as well as evaluation of the pre-training model, one-step gradient update and its low-rank approximation ($r = 8$, i.e., (Spectral-init)). Results are reported as accuracy (%) with standard deviations over 3 runs (best in **bold**). The results marked with (*) are sourced from Wang et al. (2024; 2025) under the same setting, and their hyper-parameter selection aligns with our search. The test accuracy on MNLI remains zero after one-step update thus not reported.

Method	MNLI	SST-2	CoLA	QNLI	MRPC	Avg.
Pre-train		89.79	59.03	49.28	63.48	
One-step full gradient		90.94	69.13	70.35	68.38	
$r = 8$ (low rank)		89.91	69.22	76.31	68.38	
LoRA	85.30 \pm 0.04	94.04 \pm 0.09	72.84 \pm 1.25	93.02 \pm 0.07	68.38 \pm 0.01	82.72
LoRA+*	85.81 \pm 0.09	93.85 \pm 0.24	77.53 \pm 0.20	93.14 \pm 0.03	74.43 \pm 1.39	84.95
P-LoRA	85.28 \pm 0.15	93.88 \pm 0.11	79.58 \pm 0.67	93.00 \pm 0.07	83.91 \pm 1.16	87.13
PiSSA*	85.75 \pm 0.07	94.07 \pm 0.06	74.27 \pm 0.39	93.15 \pm 0.14	76.31 \pm 0.51	84.71
LoRA-GA*	85.70 \pm 0.09	94.11 \pm 0.18	80.57 \pm 0.20	93.18 \pm 0.06	85.29 \pm 0.24	87.77
LoRA-Pro*	86.03 \pm 0.19	94.19 \pm 0.13	81.94 \pm 0.24	93.42 \pm 0.05	86.60 \pm 0.14	88.44
LoRA-One	85.89 \pm 0.08	94.53 \pm 0.13	82.04 \pm 0.22	93.37 \pm 0.02	87.83 \pm 0.37	88.73

Table 3: Performance comparison across different methods on NLG benchmarks. Results are reported as mean with standard deviations over 5 runs (higher is better).

($r = 8$)	LoRA	LoRA-GA	LoRA-One
GSM8K-D	59.26 \pm 0.99	56.44 \pm 1.15	60.44 \pm 0.17
GSM8K-CoT	53.36 \pm 0.77	46.07 \pm 1.01	55.88 \pm 0.44
MMLU	45.73 \pm 0.30	45.15 \pm 0.57	47.24 \pm 0.20
HumanEval	25.85 \pm 1.75	26.95 \pm 1.30	28.66 \pm 0.39

Table 4: The training time and memory cost of LoRA and LoRA-One from Section 6.2.

(r = 8)	Training Time (Memory)	
	LoRA	LoRA-One
MetaMathQA	6h20m (21.6GB)	6h23m (21.7GB)
Alpaca	3h22m (23.4GB)	3h25m (23.4GB)
Code-Feedback	6h24m (22.6GB)	6h26m (22.9GB)

the baselines, indicating better stability under (Spectral-init). Regarding the time and memory cost, LoRA-One takes almost the same cost as LoRA, as shown in Table 4 across all three datasets. This suggests that LoRA-One delivers its intended benefits—such as improved convergence stability or enhanced adaptability—without imposing any meaningful extra time or memory cost during fine-tuning.

6.3. Math Reasoning on Full Data and Multiple Epochs

Beyond Section 6.2, we further fine-tune LLaMA 2-7B on the complete MetaMathQA (395K) dataset for 4 epochs to access the maximum capacity of math reasoning. Here we

compare LoRA, LoRA+, LoRA-GA, and LoRA-One with rank 8. We evaluate the fine-tuned models on GSM8K with direct prompting. The learning rate and batch size are optimized for each method. More experimental details are presented in Appendix G.3.

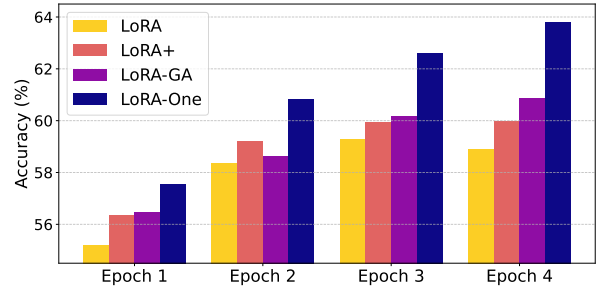


Figure 5: Accuracy comparison over epoch on GSM8K. Results are reported as mean over 2 runs (higher is better).

Fig. 5 shows that LoRA-One consistently leads other methods over epochs, suggesting that it scales more effectively with additional training and data. In contrast, LoRA-GA only shows marginal gains over LoRA and LoRA+.

7. Conclusion

This paper theoretically demonstrates how LoRA can be improved from our theoretical analysis in both linear and nonlinear models: the alignment between LoRA’s gradient update (A_t, B_t) and the singular subspace of G^\sharp , and adding preconditioners. Our theory derives the optimal initialization strategy for LoRA, clarifies some potential issues behind gradient alignment work, and bridge theory to practice with the promising performance of LoRA-One.

Impact Statement

This paper provides theoretical understanding of low-rank adapters and proposes algorithm design for parameter-efficient fine-tuning. The target of this paper is to advance the field of Machine Learning. There might be some potential societal consequences of our work, none which we feel must be specifically highlighted here.

Acknowledgment

Y. Zhang was supported by Warwick Chancellor’s International Scholarship. F. Liu was supported by Royal Society KTP R1 241011 Kan Tong Po Visiting Fellowships. Y. Chen was supported in part by National Science Foundation grants CCF-2233152. We thank Yichen Wang for coding discussions, Zulip³ for the project organization tool, and Sulis⁴ for GPU computation resources.

References

- Abbe, E., Adsera, E. B., and Misiakiewicz, T. The merged-staircase property: a necessary and nearly sufficient condition for SGD learning of sparse functions on two-layer neural networks. In *Conference on Learning Theory*, pp. 4782–4887. PMLR, 2022.
- Arous, G. B., Gheissari, R., and Jagannath, A. Online stochastic gradient descent on non-convex losses from high-dimensional inference. *Journal of Machine Learning Research*, 22(106):1–51, 2021.
- Ba, J., Erdogdu, M. A., Suzuki, T., Wang, Z., Wu, D., and Yang, G. High-dimensional Asymptotics of Feature Learning: How One Gradient Step Improves the Representation. In *Advances in Neural Information Processing Systems*, pp. 37932–37946, 2022.
- Bietti, A., Bruna, J., and Pillaud-Vivien, L. On learning gaussian multi-index models with gradient flow. *arXiv preprint arXiv:2310.19793*, 2023.
- Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J. D., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., et al. Language Models are Few-Shot Learners. In *Advances in Neural Information Processing Systems*, pp. 1877–1901, 2020.
- Brutzkus, A. and Globerson, A. Globally Optimal Gradient Descent for a ConvNet with Gaussian Inputs. In *International Conference on Machine Learning*, pp. 605–614. PMLR, 2017.
- Chen, M., Tworek, J., Jun, H., Yuan, Q., Pinto, H. P. D. O., Kaplan, J., Edwards, H., Burda, Y., Joseph, N., Brockman, G., et al. Evaluating Large Language Models Trained on Code. *arXiv preprint arXiv:2107.03374*, 2021a.
- Chen, Y., Chi, Y., Fan, J., Ma, C., et al. Spectral Methods for Data Science: A Statistical Perspective. *Foundations and Trends® in Machine Learning*, 14(5):566–806, 2021b.
- Cobbe, K., Kosaraju, V., Bavarian, M., Chen, M., Jun, H., Kaiser, L., Plappert, M., Tworek, J., Hilton, J., Nakano, R., et al. Training Verifiers to Solve Math Word Problems. *arXiv preprint arXiv:2110.14168*, 2021.
- Cui, H., Pesce, L., Dandi, Y., Krzakala, F., Lu, Y., Zdeborova, L., and Loureiro, B. Asymptotics of feature learning in two-layer networks after one gradient-step. In *International Conference on Machine Learning*, pp. 9662–9695. PMLR, 2024.
- Damian, A., Lee, J., and Soltanolkotabi, M. Neural Networks can Learn Representations with Gradient Descent. In *Conference on Learning Theory*, pp. 5413–5452. PMLR, 2022.
- Dandi, Y., Pesce, L., Cui, H., Krzakala, F., Lu, Y., and Loureiro, B. A Random Matrix Theory Perspective on the Spectrum of Learned Features and Asymptotic Generalization Capabilities. In *The 28th International Conference on Artificial Intelligence and Statistics*, 2025. URL <https://openreview.net/forum?id=3K7rKkU7Ri>.
- Dayi, A. K. and Chen, S. Gradient dynamics for low-rank fine-tuning beyond kernels. *arXiv preprint arXiv:2411.15385*, 2024.
- Ding, L., Qin, Z., Jiang, L., Zhou, J., and Zhu, Z. A Validation Approach to Over-parameterized Matrix and Image Recovery. *arXiv preprint arXiv:2209.10675*, 2022.
- Dodge, J., Ilharco, G., Schwartz, R., Farhadi, A., Hajishirzi, H., and Smith, N. Fine-Tuning Pretrained Language Models: Weight Initializations, Data Orders, and Early Stopping. *arXiv preprint arXiv:2002.06305*, 2020.
- Giampouras, P., Cai, H., and Vidal, R. Guarantees of a Preconditioned Subgradient Algorithm for Overparameterized Asymmetric Low-rank Matrix Recovery. *arXiv preprint arXiv:2410.16826*, 2024.
- Han, Z., Gao, C., Liu, J., Zhang, J., and Zhang, S. Q. Parameter-Efficient Fine-Tuning for Large Models: A Comprehensive Survey. *Transactions on Machine Learning Research*, 2024. ISSN 2835-8856. URL <https://openreview.net/forum?id=1IsCS8b6zj>.

³<https://zulip.com/>

⁴<https://warwick.ac.uk/research/rtp/sc/sulis/>

- Hayou, S., Ghosh, N., and Yu, B. LoRA+: Efficient Low Rank Adaptation of Large Models. In *International Conference on Machine Learning*, pp. 17783–17806. PMLR, 2024.
- He, K., Zhang, X., Ren, S., and Sun, J. Delving Deep into Rectifiers: Surpassing Human-Level Performance on ImageNet Classification. In *Proceedings of the IEEE international conference on computer vision*, pp. 1026–1034, 2015.
- Hendrycks, D., Burns, C., Basart, S., Zou, A., Mazeika, M., Song, D., and Steinhardt, J. Measuring Massive Multitask Language Understanding. In *International Conference on Learning Representations*, 2021.
- Horn, R. A. and Johnson, C. R. *Matrix Analysis*. Cambridge university press, 2012.
- Houlsby, N., Giurghi, A., Jastrzebski, S., Morrone, B., De Laroussilhe, Q., Gesmundo, A., Attariyan, M., and Gelly, S. Parameter-Efficient Transfer Learning for NLP. In *International Conference on Machine Learning*, pp. 2790–2799, 2019.
- Hu, E. J., Wallis, P., Allen-Zhu, Z., Li, Y., Wang, S., Wang, L., Chen, W., et al. LoRA: Low-Rank Adaptation of Large Language Models. In *International Conference on Learning Representations*, 2022.
- Jacot, A., Gabriel, F., and Hongler, C. Neural Tangent Kernel: Convergence and Generalization in Neural Networks. In *Advances in neural information processing systems*, 2018.
- Jang, U., Lee, J. D., and Ryu, E. K. LoRA Training in the NTK Regime has No Spurious Local Minima. In *International Conference on Machine Learning*, pp. 21306–21328. PMLR, 2024.
- Jia, X., Wang, H., Peng, J., Feng, X., and Meng, D. Pre-conditioning Matters: Fast Global Convergence of Non-convex Matrix Factorization via Scaled Gradient Descent. *Advances in Neural Information Processing Systems*, 36, 2024.
- Kim, J., Kim, J., and Ryu, E. K. LoRA Training Provably Converges to a Low-Rank Global Minimum or It Fails Loudly (But it Probably Won’t Fail). In *International Conference on Machine Learning*, 2025.
- Kopiczko, D. J., Blankevoort, T., and Asano, Y. M. VeRA: Vector-based Random Matrix Adaptation. In *The Twelfth International Conference on Learning Representations*, 2024.
- LeCun, Y. The mnist database of handwritten digits. <http://yann.lecun.com/exdb/mnist/>, 1998.
- Lee, J. D., Oko, K., Suzuki, T., and Wu, D. Neural network learns low-dimensional polynomials with SGD near the information-theoretic limit. *Advances in Neural Information Processing Systems*, 37:58716–58756, 2024.
- Li, B., Zhang, L., Mokhtari, A., and He, N. On the Crucial Role of Initialization for Matrix Factorization. In *The Twelfth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=YTEwJaBdh0>.
- Li, Y., Ma, T., and Zhang, H. Algorithmic Regularization in Over-parameterized Matrix Sensing and Neural Networks with Quadratic Activations. In *Conference On Learning Theory*, pp. 2–47. PMLR, 2018.
- Liu, X.-H., Du, Y., Wang, J., and Yu, Y. On the Optimization Landscape of Low Rank Adaptation Methods for Large Language Models. In *International Conference on Learning Representations*, 2025.
- Loshchilov, I. and Hutter, F. Decoupled Weight Decay Regularization. *arXiv preprint arXiv:1711.05101*, 2017.
- Ma, C., Li, Y., and Chi, Y. Beyond Procrustes: Balancing-Free Gradient Descent for Asymmetric Low-Rank Matrix Sensing. *IEEE Transactions on Signal Processing*, 69: 867–877, 2021.
- Malladi, S., Wettig, A., Yu, D., Chen, D., and Arora, S. A Kernel-Based View of Language Model Fine-Tuning. In *International Conference on Machine Learning*, pp. 23610–23641. PMLR, 2023.
- Meng, F., Wang, Z., and Zhang, M. PiSSA: Principal Singular Values and Singular Vectors Adaptation of Large Language Models. In *Advances in Neural Information Processing Systems*, pp. 121038–121072, 2024.
- Mishra, B., Apuroop, K. A., and Sepulchre, R. A riemannian geometry for low-rank matrix completion. *arXiv preprint arXiv:1211.1550*, 2012.
- Moniri, B., Lee, D., Hassani, H., and Dobriban, E. A Theory of Non-Linear Feature Learning with One Gradient Step in Two-Layer Neural Networks. In *International Conference on Machine Learning*, pp. 36106–36159. PMLR, 2024.
- Oja, E. A Simplified Neuron Model as a Principal Component Analyzer. *Journal of Mathematical Biology*, 15: 267–273, 1982.
- Oko, K., Song, Y., Suzuki, T., and Wu, D. Pretrained Transformer Efficiently Learns Low-Dimensional Target Functions In-Context. *Advances in Neural Information Processing Systems*, 37:77316–77365, 2024.

- Ponkshe, K., Singhal, R., Gorbunov, E., Tumanov, A., Horvath, S., and Vepakomma, P. Initialization using Update Approximation is a Silver Bullet for Extremely Efficient Low-Rank Fine-Tuning. *arXiv preprint arXiv:2411.19557*, 2024.
- Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M., Zhou, Y., Li, W., and Liu, P. J. Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer. *Journal of Machine Learning Research*, 21 (140):1–67, 2020.
- Soltanolkotabi, M., Stöger, D., and Xie, C. Implicit Balancing and Regularization: Generalization and Convergence Guarantees for Overparameterized Asymmetric Matrix Sensing. In *The Thirty Sixth Annual Conference on Learning Theory*, pp. 5140–5142. PMLR, 2023.
- Stöger, D. and Soltanolkotabi, M. Small random initialization is akin to spectral learning: Optimization and generalization guarantees for overparameterized low-rank matrix reconstruction. In *Advances in Neural Information Processing Systems*, pp. 23831–23843, 2021.
- Taori, R., Gulrajani, I., Zhang, T., Dubois, Y., Li, X., Guestrin, C., Liang, P., and Hashimoto, T. B. Stanford Alpaca: An Instruction-following LLaMA model. https://github.com/tatsu-lab/stanford_alpaca, 2023.
- Thoppilan, R., De Freitas, D., Hall, J., Shazeer, N., Kulshreshtha, A., Cheng, H.-T., Jin, A., Bos, T., Baker, L., Du, Y., et al. LaMDA: Language Models for Dialog Applications. *arXiv preprint arXiv:2201.08239*, 2022.
- Tong, T., Ma, C., and Chi, Y. Accelerating Ill-Conditioned Low-Rank Matrix Estimation via Scaled Gradient Descent. *Journal of Machine Learning Research*, 22(150): 1–63, 2021.
- Touvron, H., Martin, L., Stone, K., Albert, P., Almahairi, A., Babaei, Y., Bashlykov, N., Batra, S., Bhargava, P., Bhosale, S., et al. LLaMA 2: Open Foundation and Fine-Tuned Chat Models. *arXiv preprint arXiv:2307.09288*, 2023.
- Vershynin, R. Introduction to the non-asymptotic analysis of random matrices. *arXiv preprint arXiv:1011.3027*, 2010.
- Vershynin, R. *High-Dimensional Probability: An Introduction with Applications in Data Science*, volume 47. Cambridge university press, 2018.
- Wang, A., Singh, A., Michael, J., Hill, F., Levy, O., and Bowman, S. R. GLUE: A Multi-Task Benchmark and Analysis Platform for Natural Language Understanding. In *International Conference on Learning Representations*, 2019.
- Wang, S., Yu, L., and Li, J. LoRA-GA: Low-Rank Adaptation with Gradient Approximation. *Advances in Neural Information Processing Systems*, 37:54905–54931, 2024.
- Wang, Z., Liang, J., He, R., Wang, Z., and Tan, T. LoRA-Pro: Are Low-Rank Adapters Properly Optimized? In *The Twelfth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=gTwRMU3lJ5>.
- Wedin, P.-Å. Perturbation bounds in connection with singular value decomposition. *BIT Numerical Mathematics*, 12:99–111, 1972.
- Wei, J., Wang, X., Schuurmans, D., Bosma, M., Xia, F., Chi, E., Le, Q. V., Zhou, D., et al. Chain-of-Thought Prompting Elicits Reasoning in Large Language Models. *Advances in Neural Information Processing Systems*, 35: 24824–24837, 2022.
- Xiong, N., Ding, L., and Du, S. S. How Over-Parameterization Slows Down Gradient Descent in Matrix Sensing: The Curses of Symmetry and Initialization. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=xGvPKAiOhq>.
- Xu, X., Shen, Y., Chi, Y., and Ma, C. The Power of Preconditioning in Overparameterized Low-Rank Matrix Sensing. In *International Conference on Machine Learning*, pp. 38611–38654. PMLR, 2023.
- Xu, Z., Min, H., MacDonald, L. E., Luo, J., Tarmoun, S., Mallada, E., and Vidal, R. Understanding the Learning Dynamics of LoRA: A Gradient Flow Perspective on Low-Rank Adaptation in Matrix Factorization. In *International Conference on Artificial Intelligence and Statistics*, 2025.
- Yu, L., Jiang, W., Shi, H., Jincheng, Y., Liu, Z., Zhang, Y., Kwok, J., Li, Z., Weller, A., and Liu, W. MetaMath: Bootstrap Your Own Mathematical Questions for Large Language Models. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=N8N0hgNDRt>.
- Zeng, Y. and Lee, K. The Expressive Power of Low-Rank Adaptation. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=likXVjmh3E>.
- Zhang, F. and Pilanci, M. Riemannian Preconditioned LoRA for Fine-Tuning Foundation Models. In *International Conference on Machine Learning*, pp. 59641–59669. PMLR, 2024.

- Zhang, G., Fattahi, S., and Zhang, R. Y. Preconditioned Gradient Descent for Overparameterized Nonconvex Burer–Monteiro Factorization with Global Optimality Certification. *Journal of Machine Learning Research*, 24(163): 1–55, 2023.
- Zhang, J., Fattahi, S., and Zhang, R. Y. Preconditioned Gradient Descent for Over-Parameterized Nonconvex Matrix Factorization. *Advances in Neural Information Processing Systems*, 34:5985–5996, 2021.
- Zheng, T., Zhang, G., Shen, T., Liu, X., Lin, B. Y., Fu, J., Chen, W., and Yue, X. OpenCodeInterpreter: Integrating Code Generation with Execution and Refinement. In *Findings of the Association for Computational Linguistics ACL 2024*, pp. 12834–12859, 2024.
- Zhu, Z., Wu, Y., Gu, Q., and Cevher, V. Imbalance-Regularized LoRA: A Plug-and-Play Method for Improving Fine-Tuning of Foundation Models. In *NeurIPS Workshop on Adaptive Foundation Models*, 2024.

Contents

A	Symbols and Notations	15
B	Preconditioned LoRA-One	16
C	Proofs for Linear Model	17
C.1	Proofs for LoRA under Random Initialization	17
C.1.1	SVD and Schur Decomposition	18
C.1.2	Dynamics of Linear Approximation	19
C.1.3	Alignment to Negative Gradient of Full Fine-tuning	23
C.2	Gradient Descent under Spectral Initialization	27
C.3	Preconditioned Gradient Descent under Spectral Initialization	37
D	Proofs for Nonlinear Model	40
D.1	Problem Settings and Spectral Initialization	40
D.1.1	Computation of Full Population Gradients	42
D.1.2	Concentration of Empirical Gradients	46
D.2	Preconditioned Gradient Descent under Spectral Initialization	49
E	Auxiliary Results for Proofs	57
F	Detailed Comparison with LoRA-GA	58
G	Experimental Settings and Additional Results	59
G.1	Small-Scale Experiments	59
G.2	Natural Language Generation	60
G.3	Math Reasoning on Full Data and Multiple Epochs	60
G.4	Natural Language Understanding	60
G.5	Empirical Verification of Assumption 4.1	61

A. Symbols and Notations

In this section, we provide a list of the symbols and notations used in a paper.

Symbol	Dimension(s)	Definition
$\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\sigma})$	-	Multivariate normal distribution with mean vector $\boldsymbol{\mu}$ and covariance matrix $\boldsymbol{\sigma}$
$\mathcal{O}, o, \Omega, \Theta$	-	Bachmann–Landau asymptotic notation
$\ \boldsymbol{w}\ _2$	-	Euclidean norm of vector \boldsymbol{w}
$\ \mathbf{M}\ _{op}$	-	Operator norm of matrix \mathbf{M}
$\ \mathbf{M}\ _F$	-	Frobenius norm of matrix \mathbf{M}
$\langle \boldsymbol{u}, \boldsymbol{v} \rangle$	-	Dot product of vectors \boldsymbol{u} and \boldsymbol{v}
$\mathbf{M} \odot \mathbf{N}$	-	Hadamard product of matrix \mathbf{M} and \mathbf{N}
\mathbf{W}^{\natural}	$\mathbb{R}^{d \times k}$	Pre-trained weight matrix
Δ	$\mathbb{R}^{d \times k}$	Downstream feature shift matrix
$\widetilde{\mathbf{W}}^{\natural}$	$\mathbb{R}^{d \times k}$	Downstream weight matrix $\widetilde{\mathbf{W}}^{\natural} = \mathbf{W}^{\natural} + \Delta$
\mathbf{G}^{\natural}	$\mathbb{R}^{d \times k}$	The initial gradient matrix under full fine-tuning
$\mathbf{A}_t, \mathbf{B}_t$	$\mathbb{R}^{d \times r}, \mathbb{R}^{r \times k}$	Learnable low-rank adapters at step t
$\boldsymbol{w}_i^{\natural}$	\mathbb{R}^d	i^{th} column of pre-trained weight matrix \mathbf{W}^{\natural}
$\widetilde{\boldsymbol{w}}_i^{\natural}$	\mathbb{R}^d	i^{th} column of downstream weight matrix $\widetilde{\mathbf{W}}^{\natural}$
$\boldsymbol{w}_{t,i}$	\mathbb{R}^d	i^{th} column of adapted weight matrix $(\mathbf{W}^{\natural} + \mathbf{A}_t \mathbf{B}_t)$ at step t
Δ_i	\mathbb{R}^d	i^{th} column of downstream feature matrix Δ
$[\mathbf{A}_t \mathbf{B}_t]_i$	\mathbb{R}^d	i^{th} column of the product of adapters $\mathbf{A}_t \mathbf{B}_t$
$\widetilde{\mathbf{X}}$	$\mathbb{R}^{N \times d}$	Downstream data matrix
$\widetilde{\mathbf{Y}}$	$\mathbb{R}^{N \times d}$	Downstream label matrix
$\tilde{\boldsymbol{x}}_n$	\mathbb{R}^d	n^{th} downstream data point
\mathbf{M}^{-1}	-	Inverse of matrix \mathbf{M}
\mathbf{M}^{\dagger}	-	Pseudo-inverse of matrix \mathbf{M}
$\lambda_i(\mathbf{M})$	\mathbb{R}	i^{th} singular value of matrix \mathbf{M}
λ_i^*	\mathbb{R}	i^{th} singular value of downstream feature shift matrix Δ
$\kappa(\mathbf{M})$	\mathbb{R}	The condition number of matrix \mathbf{M}
κ	\mathbb{R}	The condition number of Δ : $\kappa = \lambda_{\max}^* / \lambda_{\min}^*$
κ^{\natural}	\mathbb{R}	The condition number of \mathbf{G}^{\natural} : $\kappa^{\natural} = \lambda_{\max}(\mathbf{G}^{\natural}) / \lambda_{\min}(\mathbf{G}^{\natural})$
$\mathbf{U}_m(\mathbf{M})$	-	The left singular subspace spanned by the m largest singular values of the input matrix \mathbf{M}
$\mathbf{U}_{m,\perp}(\mathbf{M})$	-	The left singular subspace orthogonal to $\mathbf{U}_m(\mathbf{M})$
$\mathbf{V}_m(\mathbf{M})$	-	The right singular subspace spanned by the m largest singular values of the input matrix \mathbf{M}
$\mathbf{V}_{m,\perp}(\mathbf{M})$	-	The right singular subspace orthogonal to $\mathbf{V}_m(\mathbf{M})$
$\mathbf{U}_{\mathbf{A}}$	-	The left singular matrix of the compact SVD of \mathbf{A}
$\mathbf{U}_{\mathbf{A},\perp}$	-	The corresponding orthogonal complement of $\mathbf{U}_{\mathbf{A}}$
$\mathbf{V}_{\mathbf{A}}$	-	The right singular matrix of the compact SVD of \mathbf{A}
$\mathbf{V}_{\mathbf{A},\perp}$	-	The corresponding orthogonal complement of $\mathbf{V}_{\mathbf{A}}$
$\sigma(\cdot)$	-	ReLU activation function
$\sigma'(\cdot)$	-	The derivative of ReLU activation function
$\nabla_{\mathbf{W}} f(\mathbf{W})$	-	The gradient matrix of function f w.r.t. input matrix \mathbf{W}
$\tilde{L}(\mathbf{A}, \mathbf{B})$	-	Loss function under LoRA fine-tuning
$L(\mathbf{W})$	-	Loss function under full fine-tuning
N	-	Number of downstream data
d	-	Input dimension of the data
k	-	Output dimension of the label
η	-	Learning rates
α	-	In theory, random init. scale of \mathbf{A}_0 . In algorithms, standard LoRA alpha.

Table 5: Essential symbols and notations in this paper.

B. Preconditioned LoRA-One

Motivated by our theory of preconditioning (see Theorem C.21 and Theorem 4.2), we propose a preconditioned variant of *LoRA-One*, termed *LoRA-One-P*. The algorithm is formally presented in Algorithm 2. *LoRA-One-P* employ the same initialization, i.e. (Spectral-init), with *LoRA-One*. For the optimizer, we use a preconditioned AdamW which introduced in (Zhang & Pilanci, 2024) instead of AdamW (Loshchilov & Hutter, 2017) in *LoRA-One*.

Algorithm 2 LoRA-One-P for a specific layer

Input: Pre-trained weight \mathbf{W}^h , batched data $\{\mathcal{D}_m\}_{m=1}^T$, sampled batch data \mathcal{B} , LoRA rank r , LoRA alpha α , loss function L , scaling parameter s , preconditioning parameter λ

Initialize:

- 1: Compute $\nabla_{\mathbf{W}} L(\mathbf{W}^h)$ given \mathcal{B}
- 2: $\mathbf{U}, \mathbf{S}, \mathbf{V} \leftarrow \text{SVD}(-\nabla_{\mathbf{W}} L(\mathbf{W}^h))$
- 3: $\mathbf{S} \leftarrow \mathbf{S}/\mathbf{S}_{[0,0]}$
- 4: $\gamma \leftarrow 1/s$
- 5: $\mathbf{A}_0 \leftarrow \sqrt{\gamma} \cdot \mathbf{U}_{[:,1:r]} \mathbf{S}_{[:,r:r]}^{1/2}$
- 6: $\mathbf{B}_0 \leftarrow \sqrt{\gamma} \cdot \mathbf{S}_{[:,r:r]}^{1/2} \mathbf{V}_{[:,1:r]}^\top$
- 7: Clear $\nabla_{\mathbf{W}} L(\mathbf{W}^h)$

Train:

- 8: **for** $t = 0, \dots, T-1$ **do**
- 9: Compute preconditioned gradients given \mathcal{D}_{t+1} :
 $\mathbf{G}_{t+1}^A \leftarrow \nabla_{\mathbf{A}} \tilde{L}(\mathbf{A}_t, \mathbf{B}_t) (\mathbf{B}_t \mathbf{B}_t^\top + \lambda \mathbf{I}_r)^{-1}$,
 $\mathbf{G}_{t+1}^B \leftarrow (\mathbf{A}_t^\top \mathbf{A}_t + \lambda \mathbf{I}_r)^{-1} \nabla_{\mathbf{B}} \tilde{L}(\mathbf{A}_t, \mathbf{B}_t)$
- 10: Update $\mathbf{A}_{t+1}, \mathbf{B}_{t+1} \leftarrow \text{AdamW}(\mathbf{G}_{t+1}^A, \mathbf{G}_{t+1}^B)$
- 11: **end for**

Return: $\mathbf{W}^h + \frac{\alpha}{\sqrt{r}} \mathbf{A}_T \mathbf{B}_T$

Next, we conduct experiments to justify that *LoRA-One-P* is more robust to sub-optimal learning rate and can achieve faster convergence under suboptimal choice than *LoRA-One*. We fine-tune the T5 base model (Raffel et al., 2020) on SST-2 dataset from GLUE (Wang et al., 2019) for one epoch. To ensure a fair comparison, we fine-tune on the grid of learning rates $\{1 \times 10^{-3}, 2 \times 10^{-4}, 1 \times 10^{-4}, 5 \times 10^{-5}, 2 \times 10^{-5}, 1 \times 10^{-5}\}$ and fix other parameters to be the same as in Appendix G.4. Additionally, we set the preconditioning parameter $\lambda = 0$ in Algorithm 2 to be consistent with Section 4. For each choice of learning rate, we run 5 different seeds for both methods and record the test accuracy every 30 steps. Then, we compute the mean and 95%-confidence interval to construct the trajectory of test accuracy during fine-tuning. The results are shown in Fig. 6.

We can observe that, *LoRA-One-P* demonstrates clear advantages over *LoRA-One* across all options of tested learning rates, and these benefits manifest in two key aspects: **robustness to mis-specified learning rates** and **faster speed of convergence**.

Under overly large learning rate (1×10^{-3}), *LoRA-One* hardly makes consistent progress and its wide confidence interval betrays unstable learning. *LoRA-One-P* makes stable improvement with a much tighter interval. This demonstrates that the preconditioners temper the volatility introduced by a too-aggressive rate, yielding both faster and more reliable gains.

When operating in the moderate range (2×10^{-4}), both methods eventually attain strong performance and share a similar trending, indicating that *LoRA-One* achieves strong performance same as *LoRA-One-P* when the learning rate is well-specified.

At the other extreme, with the learning rates chosen too small (1×10^{-4} to 1×10^{-5}), the training process of *LoRA-One* gradually appears to be slow and nearly stalls. In contrast, *LoRA-One-P* not only converges more rapidly but also maintains stable, high-quality performance even when the learning rate departs substantially from its optimal setting. This “**rescue effect**” signals that *LoRA-One-P* can tolerate sub-optimal even under-scaled updates. Its expands “safe zone” for hyperparameter tuning and reduced variance across random seeds which makes it a robust and efficient choice for real-world tasks.

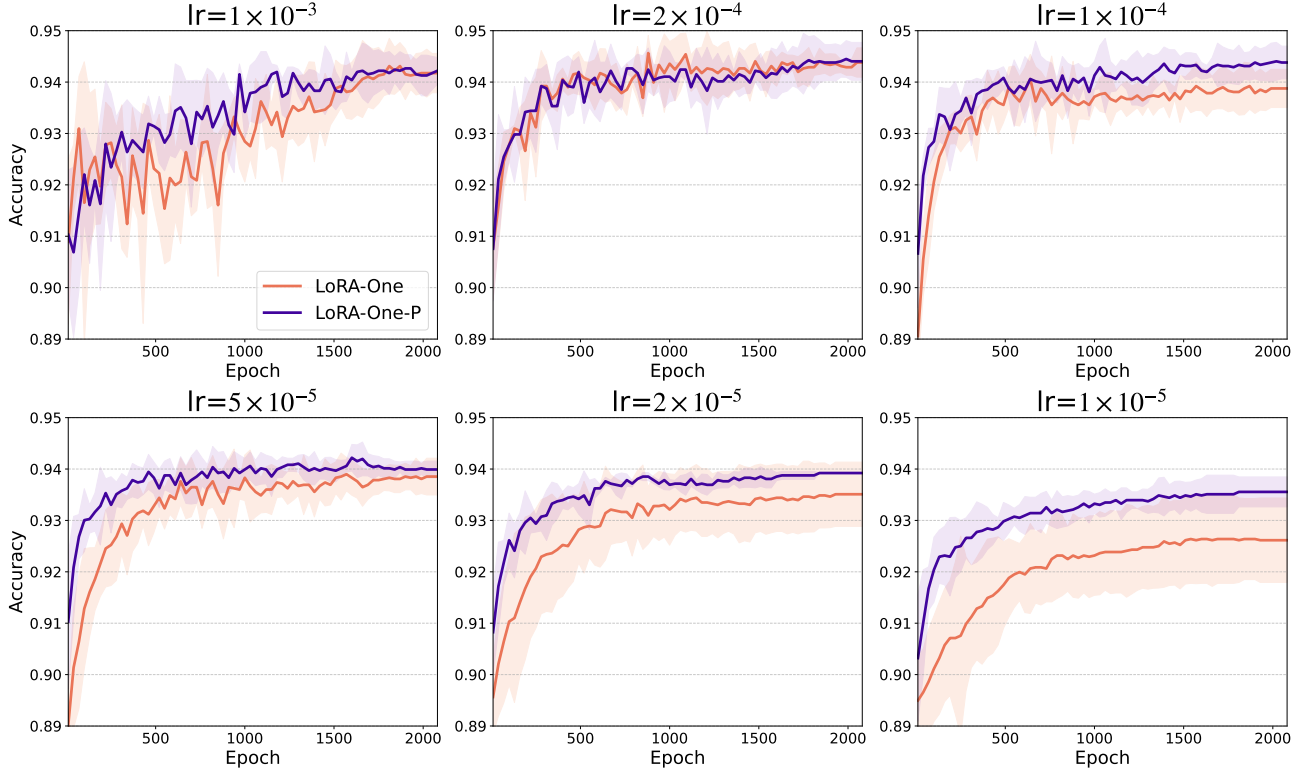


Figure 6: The trajectory of test accuracy during fine-tuning for *LoRA-One* and *LoRA-One-P* under different learning rates.

C. Proofs for Linear Model

In Appendix C.1, we deliver the proofs for alignment in Section 3.1. In Appendix C.2, we present the proofs for the main results in Section 3.2 under spectral initialization. In Appendix C.3, we give the proofs for precondition GD.

C.1. Proofs for LoRA under Random Initialization

Let $\widetilde{\mathbf{X}}$ be the fine-tuned data with $\widetilde{\mathbf{X}} \in \mathbb{R}^{N \times d}$ and the multi-output $\widetilde{\mathbf{Y}} \in \mathbb{R}^{N \times k}$. For simplicity, we define the initial residual error $\widetilde{\mathbf{Y}}_\Delta := \widetilde{\mathbf{Y}} - \widetilde{\mathbf{X}}\mathbf{W}^\natural = \widetilde{\mathbf{X}}\Delta$. Then, denote the negative gradient of Full Fine-tuning after the first step as

$$\mathbf{G}^\natural = -\nabla_{\mathbf{W}} L(\mathbf{W}^\natural) = -\frac{1}{N} \widetilde{\mathbf{X}}^\top (\widetilde{\mathbf{X}}\mathbf{W}^\natural - \widetilde{\mathbf{Y}}) = \frac{1}{N} \widetilde{\mathbf{X}}^\top \widetilde{\mathbf{Y}}_\Delta \in \mathbb{R}^{d \times k}.$$

Recall the gradient update for LoRA

$$\begin{aligned} \mathbf{A}_{t+1} &= \mathbf{A}_t - \frac{\eta}{N} \widetilde{\mathbf{X}}^\top \left(\widetilde{\mathbf{X}}(\mathbf{W}^\natural + \mathbf{A}_t \mathbf{B}_t) - \widetilde{\mathbf{Y}} \right) \mathbf{B}_t^\top, \\ \mathbf{B}_{t+1} &= \mathbf{B}_t - \frac{\eta}{N} \mathbf{A}_t^\top \widetilde{\mathbf{X}}^\top \left(\widetilde{\mathbf{X}}(\mathbf{W}^\natural + \mathbf{A}_t \mathbf{B}_t) - \widetilde{\mathbf{Y}} \right), \end{aligned}$$

we rewrite it in a compact form

$$\begin{aligned} \begin{bmatrix} \mathbf{A}_{t+1} \\ \mathbf{B}_{t+1}^\top \end{bmatrix} &= \begin{bmatrix} \mathbf{A}_t \\ \mathbf{B}_t^\top \end{bmatrix} + \underbrace{\begin{bmatrix} \mathbf{0} & \eta \mathbf{G}^\natural \\ \eta \mathbf{G}^{\natural^\top} & \mathbf{0} \end{bmatrix}}_{:=\mathbf{H}} \begin{bmatrix} \mathbf{A}_t \\ \mathbf{B}_t^\top \end{bmatrix} - \frac{\eta}{N} \begin{bmatrix} \mathbf{0} & \widetilde{\mathbf{X}}^\top \widetilde{\mathbf{X}} \mathbf{A}_t \mathbf{B}_t \\ \mathbf{B}_t^\top \mathbf{A}_t^\top \widetilde{\mathbf{X}}^\top \widetilde{\mathbf{X}} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{A}_t \\ \mathbf{B}_t^\top \end{bmatrix} \\ &= \underbrace{\begin{bmatrix} \mathbf{I}_d & \eta \mathbf{G}^\natural \\ \eta \mathbf{G}^{\natural^\top} & \mathbf{I}_k \end{bmatrix}}_{:=\mathbf{H}} \begin{bmatrix} \mathbf{A}_t \\ \mathbf{B}_t^\top \end{bmatrix} - \underbrace{\frac{\eta}{N} \begin{bmatrix} \mathbf{0} & \widetilde{\mathbf{X}}^\top \widetilde{\mathbf{X}} \mathbf{A}_t \mathbf{B}_t \\ \mathbf{B}_t^\top \mathbf{A}_t^\top \widetilde{\mathbf{X}}^\top \widetilde{\mathbf{X}} & \mathbf{0} \end{bmatrix}}_{:=\widehat{\mathbf{E}}_{t+1}} \begin{bmatrix} \mathbf{A}_t \\ \mathbf{B}_t^\top \end{bmatrix}. \end{aligned} \tag{11}$$

By defining a stack iterate

$$\mathbf{Z}_t := \begin{bmatrix} \mathbf{A}_t \\ \mathbf{B}_t^\top \end{bmatrix}, \quad \text{and} \quad \mathbf{Z}_0 := \begin{bmatrix} \mathbf{A}_0 \\ \mathbf{0} \end{bmatrix} \in \mathbb{R}^{(d+k) \times r}, \quad (12)$$

we can formulate Eq. (11) as a compact form of a nonlinear dynamical system

$$\mathbf{Z}_{t+1} = \mathbf{H} \mathbf{Z}_t - \hat{\mathbf{E}}_{t+1}, \quad (13)$$

where \mathbf{H} is a time-independent matrix corresponding to the linear part, and $\hat{\mathbf{E}}_{t+1}$ corresponds to the nonlinear part.

C.1.1. SVD AND SCHUR DECOMPOSITION

We recall the complete SVD of $\Delta \in \mathbb{R}^{d \times k}$

$$\Delta = \tilde{\mathbf{U}} \tilde{\mathbf{S}} \tilde{\mathbf{V}}^\top = [\mathbf{U} \quad \mathbf{U}_\perp] \begin{bmatrix} \mathbf{S}^* & \mathbf{0}_{r^* \times (k-r^*)} \\ \mathbf{0}_{(d-r^*) \times r^*} & \mathbf{0}_{(d-r^*) \times (k-r^*)} \end{bmatrix} \begin{bmatrix} \mathbf{V}^\top \\ \mathbf{V}_\perp^\top \end{bmatrix}, \quad \text{where } \mathbf{S}^* = \text{Diag}(\lambda_1^*, \dots, \lambda_{r^*}^*).$$

Similarly, we recall the complete SVD of \mathbf{G}^\natural as $\mathbf{G}^\natural = \tilde{\mathbf{U}}_{\mathbf{G}^\natural} \tilde{\mathbf{S}}_{\mathbf{G}^\natural} \tilde{\mathbf{V}}_{\mathbf{G}^\natural}^\top$.

We derive the Schur decomposition of \mathbf{H} under the special case $d = k$ in Lemma C.1 and then extend to $d \neq k$ in Lemma C.3 via zero padding on SVD in Lemma C.2.

Lemma C.1 (Schur Decomposition of \mathbf{H} under $d = k$). *Under assumptions in Section 2.1 for the linear setting, given $\mathbf{G}^\natural \in \mathbb{R}^{d \times k}$ in Eq. (5) and its complete SVD $\tilde{\mathbf{U}}_{\mathbf{G}^\natural} \tilde{\mathbf{S}}_{\mathbf{G}^\natural} \tilde{\mathbf{V}}_{\mathbf{G}^\natural}^\top$, if $d = k$, then the block matrix \mathbf{H} admits the following Schur decomposition*

$$\mathbf{H} = \begin{bmatrix} \mathbf{I}_d & \eta \mathbf{G}^\natural \\ \eta (\mathbf{G}^\natural)^\top & \mathbf{I}_d \end{bmatrix} = \mathbf{C} \mathbf{T} \mathbf{C}^\top,$$

where \mathbf{C} is an orthogonal matrix and \mathbf{T} is a block upper triangular matrix

$$\mathbf{C} = \frac{1}{\sqrt{2}} \begin{bmatrix} \tilde{\mathbf{U}}_{\mathbf{G}^\natural} & -\tilde{\mathbf{U}}_{\mathbf{G}^\natural} \\ \tilde{\mathbf{V}}_{\mathbf{G}^\natural} & \tilde{\mathbf{V}}_{\mathbf{G}^\natural} \end{bmatrix}, \quad \text{and} \quad \mathbf{T} = \begin{bmatrix} \mathbf{I}_d + \eta \tilde{\mathbf{S}}_{\mathbf{G}^\natural} & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_d - \eta \tilde{\mathbf{S}}_{\mathbf{G}^\natural} \end{bmatrix}.$$

Proof. We prove by verifying the claim. Starting with

$$\begin{aligned} & \begin{bmatrix} \tilde{\mathbf{U}}_{\mathbf{G}^\natural} & -\tilde{\mathbf{U}}_{\mathbf{G}^\natural} \\ \tilde{\mathbf{V}}_{\mathbf{G}^\natural} & \tilde{\mathbf{V}}_{\mathbf{G}^\natural} \end{bmatrix} \begin{bmatrix} \mathbf{I}_d + \eta \tilde{\mathbf{S}}_{\mathbf{G}^\natural} & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_d - \eta \tilde{\mathbf{S}}_{\mathbf{G}^\natural} \end{bmatrix} \\ &= \begin{bmatrix} \tilde{\mathbf{U}}_{\mathbf{G}^\natural} + \eta \tilde{\mathbf{U}}_{\mathbf{G}^\natural} \tilde{\mathbf{S}}_{\mathbf{G}^\natural} & \eta \tilde{\mathbf{U}}_{\mathbf{G}^\natural} \tilde{\mathbf{S}}_{\mathbf{G}^\natural} - \tilde{\mathbf{U}}_{\mathbf{G}^\natural} \\ \tilde{\mathbf{V}}_{\mathbf{G}^\natural} + \eta \tilde{\mathbf{V}}_{\mathbf{G}^\natural} \tilde{\mathbf{S}}_{\mathbf{G}^\natural} & \tilde{\mathbf{V}}_{\mathbf{G}^\natural} - \eta \tilde{\mathbf{V}}_{\mathbf{G}^\natural} \tilde{\mathbf{S}}_{\mathbf{G}^\natural} \end{bmatrix} =: \Xi, \end{aligned}$$

then we can verify that

$$\frac{1}{2} \times \Xi \times \begin{bmatrix} \tilde{\mathbf{U}}_{\mathbf{G}^\natural}^\top & \tilde{\mathbf{V}}_{\mathbf{G}^\natural}^\top \\ -\tilde{\mathbf{U}}_{\mathbf{G}^\natural}^\top & \tilde{\mathbf{V}}_{\mathbf{G}^\natural}^\top \end{bmatrix} = \begin{bmatrix} \mathbf{I}_d & \eta \tilde{\mathbf{U}}_{\mathbf{G}^\natural} \tilde{\mathbf{S}}_{\mathbf{G}^\natural} \tilde{\mathbf{V}}_{\mathbf{G}^\natural}^\top \\ \eta \tilde{\mathbf{V}}_{\mathbf{G}^\natural} \tilde{\mathbf{S}}_{\mathbf{G}^\natural} \tilde{\mathbf{U}}_{\mathbf{G}^\natural}^\top & \mathbf{I}_d \end{bmatrix} = \mathbf{H}.$$

Accordingly, we conclude the result. \square

Next, we consider the case of $d \neq k$.

Case 1 ($d > k$): by zero padding, \mathbf{G}^\natural and related matrices are given by

$$\underline{\mathbf{G}}^\natural = [\mathbf{G}^\natural \quad \mathbf{0}_{d \times (d-k)}], \quad \underline{\mathbf{H}} = \begin{bmatrix} \mathbf{I}_d & \eta \underline{\mathbf{G}}^\natural \\ \eta (\underline{\mathbf{G}}^\natural)^\top & \mathbf{I}_d \end{bmatrix},$$

and for any $t \geq 0$, we have the following related matrices

$$\underline{\mathbf{B}}_t = [\mathbf{B}_t \quad \mathbf{0}_{r \times (d-k)}], \quad \underline{\mathbf{Z}}_t = \begin{bmatrix} \mathbf{A}_t \\ (\underline{\mathbf{B}}_t)^\top \end{bmatrix}, \quad \underline{\mathbf{Z}}_t^{\text{lin}} = \begin{bmatrix} \mathbf{A}_t^{\text{lin}} \\ (\underline{\mathbf{B}}_t^{\text{lin}})^\top \end{bmatrix} = \underline{\mathbf{H}}^t \underline{\mathbf{Z}}_0.$$

Case 2 ($d < k$): Similarly, by zero padding, we define

$$\underline{\mathbf{G}}^{\natural} = \begin{bmatrix} \mathbf{G}^{\natural} \\ \mathbf{0}_{(k-d) \times k} \end{bmatrix}, \quad \underline{\mathbf{H}} = \begin{bmatrix} \mathbf{I}_k & \eta \underline{\mathbf{G}}^{\natural} \\ \eta (\underline{\mathbf{G}}^{\natural})^\top & \mathbf{I}_k \end{bmatrix},$$

and $\forall t \geq 0$, we define

$$\underline{\mathbf{A}}_t = \begin{bmatrix} \mathbf{A}_t \\ \mathbf{0}_{(k-d) \times r} \end{bmatrix}, \quad \underline{\mathbf{Z}}_t = \begin{bmatrix} \underline{\mathbf{A}}_t \\ (\underline{\mathbf{B}}_t)^\top \end{bmatrix}, \quad \underline{\mathbf{Z}}_t^{\text{lin}} = \begin{bmatrix} \underline{\mathbf{A}}_t^{\text{lin}} \\ (\underline{\mathbf{B}}_t^{\text{lin}})^\top \end{bmatrix} = \underline{\mathbf{H}}^t \underline{\mathbf{Z}}_0.$$

Then we have the following lemma on the SVD of $\underline{\mathbf{G}}^{\natural}$.

Lemma C.2. *If $d > k$, then we have the following SVD of $\underline{\mathbf{G}}^{\natural}$*

$$\underline{\mathbf{G}}^{\natural} = \tilde{\mathbf{U}}_{\underline{\mathbf{G}}^{\natural}} \tilde{\mathbf{S}}_{\underline{\mathbf{G}}^{\natural}} \tilde{\mathbf{V}}_{\underline{\mathbf{G}}^{\natural}}^\top,$$

where

$$\tilde{\mathbf{V}}_{\underline{\mathbf{G}}^{\natural}} = \begin{bmatrix} \tilde{\mathbf{V}}_{\mathbf{G}^{\natural}} & \mathbf{0}_{k \times (d-k)} \\ \mathbf{0}_{(d-k) \times k} & \mathbf{I}_{(d-k)} \end{bmatrix}, \quad \text{and} \quad \tilde{\mathbf{S}}_{\underline{\mathbf{G}}^{\natural}} = \begin{bmatrix} \tilde{\mathbf{S}}_{\mathbf{G}^{\natural}} & \mathbf{0}_{d \times (d-k)} \end{bmatrix}.$$

If $d < k$, then we have the following SVD of $\underline{\mathbf{G}}^{\natural}$

$$\underline{\mathbf{G}}^{\natural} = \tilde{\mathbf{U}}_{\underline{\mathbf{G}}^{\natural}} \tilde{\mathbf{S}}_{\underline{\mathbf{G}}^{\natural}} \tilde{\mathbf{V}}_{\underline{\mathbf{G}}^{\natural}}^\top,$$

where

$$\tilde{\mathbf{U}}_{\underline{\mathbf{G}}^{\natural}} = \begin{bmatrix} \tilde{\mathbf{U}}_{\mathbf{G}^{\natural}} & \mathbf{0}_{k \times (k-d)} \\ \mathbf{0}_{(k-d) \times k} & \mathbf{I}_{(k-d)} \end{bmatrix}, \quad \text{and} \quad \tilde{\mathbf{S}}_{\underline{\mathbf{G}}^{\natural}} = \begin{bmatrix} \tilde{\mathbf{S}}_{\mathbf{G}^{\natural}} \\ \mathbf{0}_{(k-d) \times k} \end{bmatrix}.$$

Proof. The block construction does not affect the original part of the SVD. It only appends zeros to the singular values and grows the corresponding orthonormal bases as partial identity matrices appropriately. \square

Now we can apply Lemma C.2 for Lemma C.1 to extend to $d \neq k$ via the following lemma. The proof is direct and we omit it here.

Lemma C.3 (Schur decomposition of $\underline{\mathbf{H}}$ under $d \neq k$). *Given the defined block matrix $\underline{\mathbf{H}} \in \mathbb{R}^{2s \times 2s}$ with $s := \max\{d, k\}$, we have the following decomposition*

$$\underline{\mathbf{H}} = \mathbf{C} \mathbf{T} \mathbf{C}^\top,$$

If $d > k$,

$$\mathbf{C} = \frac{1}{\sqrt{2}} \begin{bmatrix} \tilde{\mathbf{U}}_{\underline{\mathbf{G}}^{\natural}} & -\tilde{\mathbf{U}}_{\underline{\mathbf{G}}^{\natural}} \\ \tilde{\mathbf{V}}_{\underline{\mathbf{G}}^{\natural}} & \tilde{\mathbf{V}}_{\underline{\mathbf{G}}^{\natural}} \end{bmatrix}, \quad \mathbf{T} = \begin{bmatrix} \mathbf{I}_d + \eta \tilde{\mathbf{S}}_{\underline{\mathbf{G}}^{\natural}} & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_d - \eta \tilde{\mathbf{S}}_{\underline{\mathbf{G}}^{\natural}} \end{bmatrix}.$$

If $d < k$,

$$\mathbf{C} = \frac{1}{\sqrt{2}} \begin{bmatrix} \tilde{\mathbf{U}}_{\underline{\mathbf{G}}^{\natural}} & -\tilde{\mathbf{U}}_{\underline{\mathbf{G}}^{\natural}} \\ \tilde{\mathbf{V}}_{\underline{\mathbf{G}}^{\natural}} & \tilde{\mathbf{V}}_{\underline{\mathbf{G}}^{\natural}} \end{bmatrix}, \quad \mathbf{T} = \begin{bmatrix} \mathbf{I}_k + \eta \tilde{\mathbf{S}}_{\underline{\mathbf{G}}^{\natural}} & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_k - \eta \tilde{\mathbf{S}}_{\underline{\mathbf{G}}^{\natural}} \end{bmatrix}.$$

C.1.2. DYNAMICS OF LINEAR APPROXIMATION

The target of our proof is to demonstrate that $\hat{\mathbf{E}}_{t+1}$ does not effect the dynamics too much such that the dynamics of \mathbf{Z}_t is close to the following pseudo iterate

$$\mathbf{Z}_t^{\text{lin}} := \mathbf{H}^t \mathbf{Z}_0 =: \begin{bmatrix} \mathbf{A}_t^{\text{lin}} \\ (\mathbf{B}_t^{\text{lin}})^\top \end{bmatrix}. \quad (14)$$

The updates of the pseudo iterate follow the trajectory of Oja's Power Method (Oja, 1982). Therefore, we aim to prove that the error between the actual iterate \mathbf{Z}_t and the pseudo iterate $\mathbf{Z}_t^{\text{lin}}$ is sufficiently small, which is equivalent to that the actual iterate \mathbf{Z}_t performs a power iteration during the early steps. First, we obtain the difference between \mathbf{Z}_t and $\mathbf{Z}_t^{\text{lin}}$ by the following lemma.

Lemma C.4 (Formulation of \mathbf{E}_t). *Under assumptions in Section 2.1 for the linear setting, given the nonlinear dynamical system (13) and its linear part (14), their difference admits*

$$\mathbf{E}_t := \mathbf{Z}_t - \mathbf{Z}_t^{\text{lin}} = - \sum_{i=1}^t \mathbf{H}^{t-i} \widehat{\mathbf{E}}_i, \quad \forall t \in \mathbb{N}^+, \quad (15)$$

where $\widehat{\mathbf{E}}_i$ corresponds to the nonlinear part in Eq. (11).

Proof of Lemma C.4. We prove it by induction. Recall the formulation of the nonlinear dynamical system $\mathbf{Z}_{t+1} = \mathbf{H}\mathbf{Z}_t - \widehat{\mathbf{E}}_{t+1}$, we start with the base case $t = 1$ such that

$$\mathbf{Z}_1 = \mathbf{H}\mathbf{Z}_0 - \widehat{\mathbf{E}}_1 = \mathbf{Z}_1^{\text{lin}} - \widehat{\mathbf{E}}_1,$$

which proves the claim. Next, we assume Eq. (15) holds for $t \geq 2$, then for $t + 1$, we have

$$\begin{aligned} \mathbf{Z}_{t+1} &= \mathbf{H}\mathbf{Z}_t - \widehat{\mathbf{E}}_{t+1} \\ &= \mathbf{H} \left(\mathbf{Z}_t^{\text{lin}} - \sum_{i=1}^t \mathbf{H}^{t-i} \widehat{\mathbf{E}}_i \right) - \widehat{\mathbf{E}}_{t+1} \\ &= \mathbf{Z}_{t+1}^{\text{lin}} - \sum_{i=1}^t \mathbf{H}^{t+1-i} \widehat{\mathbf{E}}_i - \widehat{\mathbf{E}}_{t+1} \\ &= \mathbf{Z}_{t+1}^{\text{lin}} - \sum_{i=1}^{t+1} \mathbf{H}^{t+1-i} \widehat{\mathbf{E}}_i, \end{aligned}$$

which proves the claim. \square

If $\|\mathbf{E}_t\|_{op}$ is sufficiently small within a certain period, e.g., $t \leq T$, then we could approximate the early dynamics by

$$\mathbf{Z}_{t+1} := \begin{bmatrix} \mathbf{A}_{t+1} \\ \mathbf{B}_{t+1}^\top \end{bmatrix} \approx \mathbf{Z}_t^{\text{lin}} := \begin{bmatrix} \mathbf{A}_t^{\text{lin}} \\ (\mathbf{B}_t^{\text{lin}})^\top \end{bmatrix} = \begin{bmatrix} \mathbf{A}_t^{\text{lin}} \\ (\mathbf{B}_t^{\text{lin}})^\top \end{bmatrix} + \begin{bmatrix} \mathbf{0} & \eta \mathbf{G}^\natural \\ \eta \mathbf{G}^\natural^\top & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{A}_t^{\text{lin}} \\ (\mathbf{B}_t^{\text{lin}})^\top \end{bmatrix},$$

via

$$\left\| \begin{bmatrix} \mathbf{A}_t \\ \mathbf{B}_t^\top \end{bmatrix} - \begin{bmatrix} \mathbf{A}_t^{\text{lin}} \\ (\mathbf{B}_t^{\text{lin}})^\top \end{bmatrix} \right\|_{op} \leq \|\mathbf{E}_t\|_{op}.$$

In this subsection, we will bound $\|\mathbf{E}_t\|_{op}$ to show that it is actually small up to the initialization. To prove it, we first conduct the dynamical analysis of $\mathbf{Z}_t^{\text{lin}}$ via the structure of \mathbf{H} .

Part I: Dynamics of $\mathbf{Z}_t^{\text{lin}}$

With the algebra fact above, we can derive the precise spectral dynamics of $\mathbf{Z}_t^{\text{lin}}$, i.e., $\mathbf{A}_t^{\text{lin}}$ and $\mathbf{B}_t^{\text{lin}}$ separately.

Lemma C.5. *Under assumptions in Section 2.1 for the linear setting, given the pseudo iterate (14) on $\mathbf{Z}_t^{\text{lin}}$, where two components $\mathbf{A}_t^{\text{lin}}$ and $\mathbf{B}_t^{\text{lin}}$ admit the following recursion*

$$\begin{cases} \mathbf{A}_t^{\text{lin}} = \frac{1}{2} \widetilde{\mathbf{U}}_{\mathbf{G}^\natural} \left(\underbrace{\left((\mathbf{I}_d + \eta \widetilde{\mathbf{S}}_{\mathbf{G}^\natural})^t + (\mathbf{I}_d - \eta \widetilde{\mathbf{S}}_{\mathbf{G}^\natural})^t \right)}_{:= \mathbf{P}_t^A} \right) \widetilde{\mathbf{U}}_{\mathbf{G}^\natural}^\top \mathbf{A}_0, \\ (\mathbf{B}_t^{\text{lin}})^\top = \frac{1}{2} \widetilde{\mathbf{V}}_{\mathbf{G}^\natural} \left(\underbrace{\left((\mathbf{I}_d + \eta \widetilde{\mathbf{S}}_{\mathbf{G}^\natural})^t - (\mathbf{I}_d - \eta \widetilde{\mathbf{S}}_{\mathbf{G}^\natural})^t \right)}_{:= \mathbf{P}_t^B} \right) \widetilde{\mathbf{U}}_{\mathbf{G}^\natural}^\top \mathbf{A}_0. \end{cases}$$

Furthermore, if $\widetilde{\mathbf{X}}^\top \widetilde{\mathbf{X}}$ is non-singular, \mathbf{P}_t^A is a full rank matrix and singular values are 1 after the r^* -th order. \mathbf{P}_t^B is a rank- r^* matrix.

Proof. We start with the special case $d = k$ and then discuss the case of $d \neq k$. For the case of $d = k$, we have

$$\mathbf{Z}_t^{\text{lin}} = \mathbf{H}^t \mathbf{Z}_0 = (\mathbf{C} \mathbf{T} \mathbf{C}^\top)^t \mathbf{Z}_0 = \mathbf{C} \mathbf{T}^t \mathbf{C}^\top \mathbf{Z}_0,$$

where the last equality follows from the fact that \mathbf{C} is an orthogonal matrix. Next, we compute \mathbf{T}^t

$$\mathbf{T}^t = \begin{bmatrix} \left(\mathbf{I}_d + \eta \tilde{\mathbf{S}}_{\mathbf{G}^\natural} \right)^t & \mathbf{0} \\ \mathbf{0}_{d \times d} & \left(\mathbf{I}_d - \eta \tilde{\mathbf{S}}_{\mathbf{G}^\natural} \right)^t \end{bmatrix}. \quad (16)$$

Then, we can derive the following recursion

$$\begin{aligned} \mathbf{Z}_t^{\text{lin}} &= \mathbf{H}^t \mathbf{Z}_0 \\ &= \frac{1}{\sqrt{2}} \begin{bmatrix} \tilde{\mathbf{U}}_{\mathbf{G}^\natural} & -\tilde{\mathbf{U}}_{\mathbf{G}^\natural} \\ \tilde{\mathbf{V}}_{\mathbf{G}^\natural} & \tilde{\mathbf{V}}_{\mathbf{G}^\natural} \end{bmatrix} \begin{bmatrix} \left(\mathbf{I}_d + \eta \tilde{\mathbf{S}}_{\mathbf{G}^\natural} \right)^t & \mathbf{0} \\ \mathbf{0}_{d \times d} & \left(\mathbf{I}_d - \eta \tilde{\mathbf{S}}_{\mathbf{G}^\natural} \right)^t \end{bmatrix} \times \mathbf{C}^\top \mathbf{Z}_0 \\ &= \begin{bmatrix} \tilde{\mathbf{U}}_{\mathbf{G}^\natural} \left(\mathbf{I}_d + \eta \tilde{\mathbf{S}}_{\mathbf{G}^\natural} \right)^t & -\tilde{\mathbf{U}}_{\mathbf{G}^\natural} \left(\mathbf{I}_d - \eta \tilde{\mathbf{S}}_{\mathbf{G}^\natural} \right)^t \\ \tilde{\mathbf{V}}_{\mathbf{G}^\natural} \left(\mathbf{I}_d + \eta \tilde{\mathbf{S}}_{\mathbf{G}^\natural} \right)^t & \tilde{\mathbf{V}}_{\mathbf{G}^\natural} \left(\mathbf{I}_d - \eta \tilde{\mathbf{S}}_{\mathbf{G}^\natural} \right)^t \end{bmatrix} \times \frac{\mathbf{C}^\top \mathbf{Z}_0}{\sqrt{2}} \\ &= \begin{bmatrix} \frac{1}{2} \tilde{\mathbf{U}}_{\mathbf{G}^\natural} \left(\left(\mathbf{I}_d + \eta \tilde{\mathbf{S}}_{\mathbf{G}^\natural} \right)^t + \left(\mathbf{I}_d - \eta \tilde{\mathbf{S}}_{\mathbf{G}^\natural} \right)^t \right) \tilde{\mathbf{U}}_{\mathbf{G}^\natural}^\top & * \\ \frac{1}{2} \tilde{\mathbf{V}}_{\mathbf{G}^\natural} \left(\left(\mathbf{I}_d + \eta \tilde{\mathbf{S}}_{\mathbf{G}^\natural} \right)^t - \left(\mathbf{I}_d - \eta \tilde{\mathbf{S}}_{\mathbf{G}^\natural} \right)^t \right) \tilde{\mathbf{U}}_{\mathbf{G}^\natural}^\top & * \end{bmatrix} \begin{bmatrix} \mathbf{A}_0 \\ \mathbf{0} \end{bmatrix} \\ &= \begin{bmatrix} \frac{1}{2} \tilde{\mathbf{U}}_{\mathbf{G}^\natural} \left(\left(\mathbf{I}_d + \eta \tilde{\mathbf{S}}_{\mathbf{G}^\natural} \right)^t + \left(\mathbf{I}_d - \eta \tilde{\mathbf{S}}_{\mathbf{G}^\natural} \right)^t \right) \tilde{\mathbf{U}}_{\mathbf{G}^\natural}^\top \mathbf{A}_0 \\ \frac{1}{2} \tilde{\mathbf{V}}_{\mathbf{G}^\natural} \left(\left(\mathbf{I}_d + \eta \tilde{\mathbf{S}}_{\mathbf{G}^\natural} \right)^t - \left(\mathbf{I}_d - \eta \tilde{\mathbf{S}}_{\mathbf{G}^\natural} \right)^t \right) \tilde{\mathbf{U}}_{\mathbf{G}^\natural}^\top \mathbf{A}_0 \end{bmatrix}. \end{aligned}$$

Next, we extend the results above to $d \neq k$. Here we take $d > k$,

$$\begin{aligned} \mathbf{B}_t^{\text{lin}} &= \frac{1}{2} \tilde{\mathbf{V}}_{\mathbf{G}^\natural} \left(\left(\mathbf{I}_d + \eta \tilde{\mathbf{S}}_{\mathbf{G}^\natural} \right)^t - \left(\mathbf{I}_d - \eta \tilde{\mathbf{S}}_{\mathbf{G}^\natural} \right)^t \right) \tilde{\mathbf{U}}_{\mathbf{G}^\natural}^\top \mathbf{A}_0 \\ &= \begin{bmatrix} \frac{1}{2} \tilde{\mathbf{V}}_{\mathbf{G}^\natural} \left(\left(\mathbf{I}_d + \eta \tilde{\mathbf{S}}_{\mathbf{G}^\natural} \right)^t - \left(\mathbf{I}_d - \eta \tilde{\mathbf{S}}_{\mathbf{G}^\natural} \right)^t \right) \tilde{\mathbf{U}}_{\mathbf{G}^\natural}^\top \mathbf{A}_0 & \mathbf{0}_{r \times (d-k)} \end{bmatrix}, \end{aligned}$$

which proves the claim. Lastly, we take $d < k$,

$$\begin{aligned} \mathbf{A}_t^{\text{lin}} &= \frac{1}{2} \tilde{\mathbf{U}}_{\mathbf{G}^\natural} \left(\left(\mathbf{I}_k + \eta \tilde{\mathbf{S}}_{\mathbf{G}^\natural} \right)^t + \left(\mathbf{I}_k - \eta \tilde{\mathbf{S}}_{\mathbf{G}^\natural} \right)^t \right) \tilde{\mathbf{U}}_{\mathbf{G}^\natural}^\top \mathbf{A}_0 \\ &\quad \begin{bmatrix} \frac{1}{2} \tilde{\mathbf{U}}_{\mathbf{G}^\natural} \left(\left(\mathbf{I}_d + \eta \tilde{\mathbf{S}}_{\mathbf{G}^\natural} \right)^t + \left(\mathbf{I}_d - \eta \tilde{\mathbf{S}}_{\mathbf{G}^\natural} \right)^t \right) \tilde{\mathbf{U}}_{\mathbf{G}^\natural}^\top \mathbf{A}_0 \\ \mathbf{0}_{(k-d) \times r} \end{bmatrix}, \end{aligned}$$

which completes the proof.

Besides, we discuss about some properties of $\mathbf{P}_t^{\mathbf{A}}$ and $\mathbf{P}_t^{\mathbf{B}}$. Recall $\text{Rank}(\mathbf{G}^\natural) = \text{Rank}(\Delta) = r^*$, then we have

$$\lambda_{r^*+i}(\mathbf{P}_t^{\mathbf{A}}) = \frac{1}{2} \lambda_{r^*+i} \left(\left(\mathbf{I}_d + \eta \tilde{\mathbf{S}}_{\mathbf{G}^\natural} \right)^t + \left(\mathbf{I}_d - \eta \tilde{\mathbf{S}}_{\mathbf{G}^\natural} \right)^t \right) = 1, \quad \forall 1 \leq i \leq (d - r^*).$$

That means $\mathbf{P}_t^{\mathbf{A}} \in \mathbb{R}^{d \times d}$ is a full rank matrix and the singular values are 1 after the r^* -th order. However $\mathbf{P}_t^{\mathbf{B}} \in \mathbb{R}^{k \times k}$ is a rank- r^* matrix. \square

Part II: Control $\|\mathbf{E}_t\|_{op}$

Based on the above results, we are ready to prove that $\|\mathbf{E}_t\|_{op}$ is small.

Lemma C.6. Under assumptions in Section 2.1 for the linear setting, with LoRA initialization (LoRA-init), given $\|\mathbf{A}_0\|_{op}$ and \mathbf{G}^\natural in Eq. (5) and its largest singular value $\lambda_1(\mathbf{G}^\natural)$, consider the following time period

$$t \leq t^* := \frac{\ln \left(\frac{\lambda_1(\mathbf{G}^\natural)}{3\|\mathbf{A}_0\|_{op}^2} \right)}{3 \ln(1 + \eta\lambda_1(\mathbf{G}^\natural))},$$

then the following statement holds with probability at least $1 - 2C \exp(-N)$ for a universal constant C over random Gaussian data

$$\|\mathbf{E}_t\|_{op} \leq \|\mathbf{A}_0\|_{op}. \quad (17)$$

Remark: By choosing proper random initialization variance over \mathbf{A}_0 , we can ensure $t^* > 1$ to avoid vacuous upper bound.

Proof. We will prove by induction. Starting from $t = 0$, this is trivially true since $\mathbf{Z}_0 = \mathbf{Z}_0^{\text{lin}}$. Next, we assume Eq. (17) holds for $t - 1$ with $t \geq 1$ and prove $\|\mathbf{E}_t\|_{op} \leq \|\mathbf{A}_0\|_{op}$. To deliver the proof, denote $a_0 := \|\mathbf{A}_0\|_{op}$, from Lemma C.5, we know that

$$\|\mathbf{A}_{t-1}^{\text{lin}}\|_{op} \leq (1 + \eta\lambda_1(\mathbf{G}^\natural))^{t-1} a_0, \quad \|\mathbf{B}_{t-1}^{\text{lin}}\|_{op} \leq \frac{1}{2} (1 + \eta\lambda_1(\mathbf{G}^\natural))^{t-1} a_0. \quad (18)$$

Besides, since $(\mathbf{A}_t - \mathbf{A}_t^{\text{lin}})$ and $(\mathbf{B}_t - \mathbf{B}_t^{\text{lin}})$ are the sub-matrices of the error term \mathbf{E}_t , our condition $\|\mathbf{E}_{t-1}\|_{op} \leq \|\mathbf{A}_0\|_{op}$ we have

$$\begin{cases} \|\mathbf{A}_{t-1} - \mathbf{A}_{t-1}^{\text{lin}}\|_{op} \leq \|\mathbf{E}_{t-1}\|_{op}, \\ \|\mathbf{B}_{t-1} - \mathbf{B}_{t-1}^{\text{lin}}\|_{op} \leq \|\mathbf{E}_{t-1}\|_{op}. \end{cases} \quad (19)$$

It implies that

$$\|\mathbf{A}_{t-1}\|_{op} \leq (1 + \eta\lambda_1(\mathbf{G}^\natural))^{t-1} a_0 + \|\mathbf{E}_{t-1}\|_{op}, \quad \|\mathbf{B}_{t-1}\|_{op} \leq \frac{1}{2} (1 + \eta\lambda_1(\mathbf{G}^\natural))^{t-1} a_0 + \|\mathbf{E}_{t-1}\|_{op}.$$

Besides, according to covariance matrix estimation in the operator norm in Lemma E.1, with probability at least $1 - 2C \exp(-N\epsilon^2)$ for a universal constant $C > 0$, we have (taking $\epsilon = 1$)

$$\left\| \frac{1}{N} \widetilde{\mathbf{X}}^\top \widetilde{\mathbf{X}} - \mathbf{I}_d \right\|_{op} \leq \epsilon = 1. \quad (20)$$

Accordingly, with probability at least $1 - 2C \exp(-N)$, $\|\widehat{\mathbf{E}}_t\|_{op}$ can be upper bounded by

$$\begin{aligned} \|\widehat{\mathbf{E}}_t\|_{op} &\leq \eta \left\| \frac{1}{N} \widetilde{\mathbf{X}}^\top \widetilde{\mathbf{X}} \mathbf{A}_{t-1} \mathbf{B}_{t-1} \mathbf{B}_{t-1}^\top \right\|_{op} + \eta \left\| \mathbf{B}_{t-1}^\top \mathbf{A}_{t-1}^\top \frac{1}{N} \widetilde{\mathbf{X}}^\top \widetilde{\mathbf{X}} \mathbf{A}_{t-1} \right\|_{op} \\ &\leq \eta(1 + \epsilon) \|\mathbf{A}_{t-1}\|_{op} \|\mathbf{B}_{t-1}\|_{op}^2 + \eta(1 + \epsilon) \|\mathbf{A}_{t-1}\|_{op}^2 \|\mathbf{B}_{t-1}\|_{op} \quad [\text{using Eq. (20)}] \\ &\leq (1 + \epsilon) \eta \|\mathbf{A}_{t-1}\|_{op} \|\mathbf{B}_{t-1}\|_{op} (\|\mathbf{B}_{t-1}\|_{op} + \|\mathbf{A}_{t-1}\|_{op}) \\ &\leq (1 + \epsilon) \eta (\|\mathbf{A}_{t-1}^{\text{lin}}\|_{op} + \|\mathbf{E}_{t-1}\|_{op}) (\|\mathbf{B}_{t-1}^{\text{lin}}\|_{op} + \|\mathbf{E}_{t-1}\|_{op}) \times (\|\mathbf{B}_{t-1}^{\text{lin}}\|_{op} + \|\mathbf{A}_{t-1}^{\text{lin}}\|_{op} + 2\|\mathbf{E}_{t-1}\|_{op}) \quad [\text{using Eq. (19)}] \end{aligned}$$

Accordingly, using the upper bound of $\|\mathbf{A}_{t-1}^{\text{lin}}\|_{op}$ and $\|\mathbf{B}_{t-1}^{\text{lin}}\|_{op}$ in Eq. (18), we have

$$\begin{aligned} \|\widehat{\mathbf{E}}_t\|_{op} &\leq (1 + \epsilon) \eta \left((1 + \eta\lambda_1(\mathbf{G}^\natural))^{t-1} a_0 + \|\mathbf{E}_{t-1}\|_{op} \right) \left(\frac{1}{2} (1 + \eta\lambda_1(\mathbf{G}^\natural))^{t-1} a_0 + \|\mathbf{E}_{t-1}\|_{op} \right) \times \\ &\quad \left(\frac{3}{2} (1 + \eta\lambda_1(\mathbf{G}^\natural))^{t-1} a_0 + 2\|\mathbf{E}_{t-1}\|_{op} \right) \\ &\leq 2(1 + \epsilon) \eta \left((1 + \eta\lambda_1(\mathbf{G}^\natural))^{3t-3} a_0^3 + \|\mathbf{E}_{t-1}\|_{op}^3 \right) \\ &\leq 6\eta (1 + \eta\lambda_1(\mathbf{G}^\natural))^{3t-3} a_0^3. \quad [\text{from our inductive hypothesis}] \end{aligned}$$

Then, by Lemma C.4, we can conclude that

$$\begin{aligned}
 \|\mathbf{E}_t\|_{op} &= \left\| \sum_{i=1}^t \mathbf{H}^{t-i} \widehat{\mathbf{E}}_i \right\|_{op} \leq \sum_{i=1}^t \|\mathbf{H}\|_{op}^{t-i} \|\widehat{\mathbf{E}}_i\|_{op} \\
 &\leq 6\eta a_0^3 \times \sum_{i=1}^t (1 + \eta\lambda_1(\mathbf{G}^\natural))^{t+2i-3} \quad [\text{using Lemma C.1}] \\
 &= 6\eta a_0^3 \times (1 + \eta\lambda_1(\mathbf{G}^\natural))^{t-1} \sum_{i=1}^t (1 + \eta\lambda_1(\mathbf{G}^\natural))^{2i-2} \\
 &= 6\eta a_0^3 \times (1 + \eta\lambda_1(\mathbf{G}^\natural))^{t-1} \frac{(1 + \eta\lambda_1(\mathbf{G}^\natural))^{2t} - 1}{(1 + \eta\lambda_1(\mathbf{G}^\natural))^2 - 1} \quad [\text{geometric series}] \\
 &\leq 6\eta a_0^3 \times (1 + \eta\lambda_1(\mathbf{G}^\natural))^{t-1} \frac{(1 + \eta\lambda_1(\mathbf{G}^\natural))^{2t+1}}{2\eta\lambda_1(\mathbf{G}^\natural)} \\
 &\leq 3 (1 + \eta\lambda_1(\mathbf{G}^\natural))^{3t} \frac{a_0^3}{\lambda_1(\mathbf{G}^\natural)}. \tag{21}
 \end{aligned}$$

Accordingly, when $t \leq t^* := \frac{\ln\left(\frac{\lambda_1(\mathbf{G}^\natural)}{3\|\mathbf{A}_0\|_{op}^2}\right)}{3\ln(1+\eta\lambda_1(\mathbf{G}^\natural))}$, we have

$$\|\mathbf{E}_t\|_{op} \leq \|\mathbf{A}_0\|_{op},$$

which proves the claim. \square

C.1.3. ALIGNMENT TO NEGATIVE GRADIENT OF FULL FINE-TUNING

Now we can apply Lemma C.6 to obtain

$$\|\mathbf{A}_t - \mathbf{A}_t^{\text{lin}}\|_{op} \leq \|\mathbf{A}_0\|_{op}.$$

Recall Lemma C.5, we can observe that the dynamic of $\mathbf{A}_t^{\text{lin}}$ also follows an Oja's Power Method (Oja, 1982), which aligns $\mathbf{A}_t^{\text{lin}}$'s left singular subspace to the left subspace of the initial negative gradient step \mathbf{G}^\natural of full fine-tuning. We anticipate that $\lambda_{r^*}(\mathbf{A}_t) \gg \lambda_{r^*+1}(\mathbf{A}_t)$ for sufficiently large t . Furthermore, if $\|\mathbf{E}_t\|_{op}$ remains small, then the top- r^* left singular subspace of \mathbf{A}_t can closely align to \mathbf{G}^\natural 's. To prove this alignment, we modify Stöger & Soltanolkotabi (2021, Lemma 8.3) to obtain the following results.

Lemma C.7. *Under assumptions in Section 2.1 for the linear setting, recall $\mathbf{P}_t^{\mathbf{A}} := \frac{1}{2}\tilde{\mathbf{U}}_{\mathbf{G}^\natural} \left((\mathbf{I}_d + \eta\tilde{\mathbf{S}}_{\mathbf{G}^\natural})^t + (\mathbf{I}_d - \eta\tilde{\mathbf{S}}_{\mathbf{G}^\natural})^t \right) \tilde{\mathbf{U}}_{\mathbf{G}^\natural}^\top$ as $\mathbb{R}^{d \times d}$ -valued symmetric matrix in Lemma C.5, we assume that*

$$\lambda_{r^*+1}(\mathbf{P}_t^{\mathbf{A}})\|\mathbf{A}_0\|_{op} + \|\mathbf{E}_t\|_{op} < \lambda_{r^*}(\mathbf{P}_t^{\mathbf{A}})\lambda_{\min}(\mathbf{U}_{r^*}^\top(\mathbf{P}_t^{\mathbf{A}})\mathbf{A}_0),$$

that can be satisfied under certain conditions (discussed later). Then the following three inequalities hold:

$$\lambda_{r^*}(\mathbf{P}_t^{\mathbf{A}}\mathbf{A}_0 + \mathbf{E}_t) \geq \lambda_{r^*}(\mathbf{P}_t^{\mathbf{A}})\lambda_{\min}(\mathbf{U}_{r^*}^\top(\mathbf{P}_t^{\mathbf{A}})\mathbf{A}_0) - \|\mathbf{E}_t\|_{op}, \tag{22}$$

$$\lambda_{r^*+1}(\mathbf{P}_t^{\mathbf{A}}\mathbf{A}_0 + \mathbf{E}_t) \leq \lambda_{r^*+1}(\mathbf{P}_t^{\mathbf{A}})\|\mathbf{A}_0\|_{op} + \|\mathbf{E}_t\|_{op}, \tag{23}$$

$$\|\mathbf{U}_{r^*,\perp}^\top(\mathbf{P}_t^{\mathbf{A}})\mathbf{U}_{r^*}(\mathbf{P}_t^{\mathbf{A}}\mathbf{A}_0 + \mathbf{E}_t)\|_{op} \leq \frac{\lambda_{r^*+1}(\mathbf{P}_t^{\mathbf{A}})\|\mathbf{A}_0\|_{op} + \|\mathbf{E}_t\|_{op}}{\lambda_{r^*}(\mathbf{P}_t^{\mathbf{A}})\lambda_{\min}(\mathbf{U}_{r^*}^\top(\mathbf{P}_t^{\mathbf{A}})\mathbf{A}_0) - \lambda_{r^*+1}(\mathbf{P}_t^{\mathbf{A}})\|\mathbf{A}_0\|_{op} - \|\mathbf{E}_t\|_{op}}, \tag{24}$$

where $\mathbf{U}_k(\mathbf{M})$ denotes the left singular subspace spanned by the k largest singular values of the input matrix \mathbf{M} and $\mathbf{U}_{k,\perp}(\mathbf{M})$ denotes the left singular subspace orthogonal to $\mathbf{U}_k(\mathbf{M})$.

This lemma can help us derive the principle angle of the left singular subspace between $\mathbf{A}_t^{\text{lin}}$ and \mathbf{A}_t . Note that the assumption comes from the necessary condition of Wedin's $\sin \theta$ theorem (Wedin, 1972). In the next lemma, we aim to derive the time threshold which can fulfill this assumption.

Lemma C.8. Under assumptions in Section 2.1 for the linear setting, given $\|\mathbf{A}_0\|_{op}$, for any $\theta \in (0, 1)$, taking

$$t \leq \frac{\ln \left(\frac{8\|\mathbf{A}_0\|_{op}}{\theta \lambda_{\min}(\mathbf{U}_{r^*}^\top (\mathbf{P}_t^{\mathbf{A}}) \mathbf{A}_0)} \right)}{\ln(1 + \eta \lambda_{r^*}(\mathbf{G}^{\natural}))},$$

then Eq. (24) holds with probability at least $1 - 2C \exp(-N)$ for a universal constant C over random Gaussian data, i.e.

$$\|\mathbf{U}_{r^*, \perp}^\top (\mathbf{P}_t^{\mathbf{A}}) \mathbf{U}_{r^*} (\mathbf{P}_t^{\mathbf{A}} \mathbf{A}_0 + \mathbf{E}_t)\|_{op} \leq \theta.$$

Remark: To ensure that the θ -alignment phase still falls into the early phase in Lemma C.6 for $\|\mathbf{E}_t\|_{op} \leq \|\mathbf{A}_0\|_{op}$, we need to choose proper initialization for \mathbf{A}_0 . We will detail this in Theorem 3.2 later.

Proof. First, $\lambda_{r^*}(\mathbf{P}_t^{\mathbf{A}})$ in Lemma C.5 can be lower bounded by

$$\begin{aligned} \lambda_{r^*}(\mathbf{P}_t^{\mathbf{A}}) &= \frac{1}{2} \lambda_{r^*} \left((\mathbf{I}_d + \eta \tilde{\mathbf{S}}_{\mathbf{G}^{\natural}})^t + (\mathbf{I}_d - \eta \tilde{\mathbf{S}}_{\mathbf{G}^{\natural}})^t \right) \\ &\geq \frac{1}{2} \lambda_{r^*} \left((\mathbf{I}_d + \eta \tilde{\mathbf{S}}_{\mathbf{G}^{\natural}})^t \right) \\ &= \frac{1}{2} (1 + \eta \lambda_{r^*}(\mathbf{G}^{\natural}))^t. \end{aligned} \quad (25)$$

Recall Lemma C.5, we have $\lambda_{r^*+1}(\mathbf{P}_t^{\mathbf{A}}) = 1$ and Lemma C.6 with $\|\mathbf{E}_t\|_{op} \leq \|\mathbf{A}_0\|_{op}$, we define the following threshold γ and upper bound it

$$\begin{aligned} \gamma &:= \frac{\lambda_{r^*+1}(\mathbf{P}_t^{\mathbf{A}}) \|\mathbf{A}_0\|_{op} + \|\mathbf{E}_t\|_{op}}{\lambda_{r^*}(\mathbf{P}_t^{\mathbf{A}}) \lambda_{\min}(\mathbf{U}_{r^*}^\top (\mathbf{P}_t^{\mathbf{A}}) \mathbf{A}_0)} \\ &\leq \frac{2\|\mathbf{A}_0\|_{op}}{\frac{1}{2} (1 + \eta \lambda_{r^*}(\mathbf{G}^{\natural}))^t \lambda_{\min}(\mathbf{U}_{r^*}^\top (\mathbf{P}_t^{\mathbf{A}}) \mathbf{A}_0)} \quad [\text{using Lemma C.5, C.6}] \\ &= \exp(-\ln(1 + \eta \lambda_{r^*}(\mathbf{G}^{\natural})) \cdot t) \cdot \frac{4\|\mathbf{A}_0\|_{op}}{\lambda_{\min}(\mathbf{U}_{r^*}^\top (\mathbf{P}_t^{\mathbf{A}}) \mathbf{A}_0)}. \end{aligned} \quad (26)$$

Set $\theta \in (0, 1)$, let $\text{Eq.}(26) \leq \frac{\theta}{2}$, then we have that

$$\|\mathbf{U}_{r^*, \perp}^\top (\mathbf{P}_t^{\mathbf{A}}) \mathbf{U}_{r^*} (\mathbf{P}_t^{\mathbf{A}} \mathbf{A}_0 + \mathbf{E}_t)\|_{op} \leq \theta.$$

The time t to achieve this angle θ can be upper bounded by

$$\exp(-\ln(1 + \eta \lambda_{r^*}(\mathbf{G}^{\natural})) \cdot t) \cdot \frac{4\|\mathbf{A}_0\|_{op}}{\lambda_{\min}(\mathbf{U}_{r^*}^\top (\mathbf{P}_t^{\mathbf{A}}) \mathbf{A}_0)} \leq \frac{\theta}{2},$$

which implies that

$$t \leq \frac{\ln \left(\frac{8\|\mathbf{A}_0\|_{op}}{\theta \lambda_{\min}(\mathbf{U}_{r^*}^\top (\mathbf{P}_t^{\mathbf{A}}) \mathbf{A}_0)} \right)}{\ln(1 + \eta \lambda_{r^*}(\mathbf{G}^{\natural}))}.$$

Finally we conclude the proof. \square

Theorem C.9. [Full version of Theorem 3.2] Under assumptions in Section 2.1 for the linear setting, recall \mathbf{G}^{\natural} defined in Eq. (5) with its condition number κ^{\natural} , we consider random Gaussian initialization $\mathbf{A}_0 \in \mathbb{R}^{d \times r}$ with $[\mathbf{A}_0]_{ij} \sim \mathcal{N}(0, \alpha^2)$ in (LoRA-init), for any $\theta \in (0, 1)$, let $\xi = o(1)$ be chosen such that

$$\alpha \leq \begin{cases} \left(\frac{\theta \xi}{24r\sqrt{d}} \right)^{\frac{3\kappa^{\natural}}{2}} \sqrt{\frac{\lambda_1(\mathbf{G}^{\natural})}{27d}} & \text{if } r^* \leq r < 2r^*, \\ \left(\frac{\theta}{24\sqrt{d}} \right)^{\frac{3\kappa^{\natural}}{2}} \sqrt{\frac{\lambda_1(\mathbf{G}^{\natural})}{27d}} & \text{if } r \geq 2r^*. \end{cases}$$

Then if we run gradient descent for t^* steps with

$$t^* \lesssim \begin{cases} \frac{\ln\left(\frac{24r\sqrt{d}}{\theta\xi}\right)}{\ln(1+\eta\lambda_{r^*}(\mathbf{G}^\natural))} & \text{if } r^* \leq r < 2r^*, \\ \frac{\ln\left(\frac{24\sqrt{d}}{\theta}\right)}{\ln(1+\eta\lambda_{r^*}(\mathbf{G}^\natural))} & \text{if } r \geq 2r^*, \end{cases}$$

we have the following alignment on the left singular subspace between \mathbf{G}^\natural and \mathbf{A}_{t^*}

$$\begin{aligned} & \|\mathbf{U}_{r^*,\perp}^\top(\mathbf{G}^\natural)\mathbf{U}_{r^*}(\mathbf{A}_{t^*})\|_{op} \lesssim \theta, \\ & \text{with probability at least } \begin{cases} 1 - C_1 \exp(-d) - (C_2\xi)^{r-r^*+1} - C_3 \exp(-r) - C \exp(-N) & \text{if } r^* \leq r < 2r^*, \\ 1 - C_4 \exp(-d) - C_5 \exp(-r) - C \exp(-N) & \text{if } r \geq 2r^*, \end{cases} \end{aligned}$$

for some positive constants $C, C_1, C_2, C_3, C_4, C_5$. Here $\mathbf{U}_{r^*}(\mathbf{A}_{t^*})$ denotes the left singular subspace spanned by the r^* largest singular values of \mathbf{A}_{t^*} and $\mathbf{U}_{r^*,\perp}(\mathbf{M})$ denotes the left singular subspace orthogonal to $\mathbf{U}_{r^*}(\mathbf{M})$. Note that we can select any pair of stepsizes (η, η) that satisfies the conditions $t^* > 1$, $\eta \geq \eta$, and $\zeta(\eta, \eta) = \Theta(1)$.

Proof. For ease of description, we denote $\mathbf{A}_0 := \alpha \mathbf{T} \in \mathbb{R}^{d \times r}$ where \mathbf{T} is a standard random Gaussian matrix with zero-mean and unit variance. Here we aim to choose a proper α to ensure that θ -alignment phase in Lemma C.8 still falls into the early phase in Lemma C.6, i.e.

$$\begin{aligned} & \frac{\ln\left(\frac{8\|\mathbf{A}_0\|_{op}}{\theta\lambda_{\min}(\mathbf{U}_{r^*}^\top(\mathbf{P}_t^{\mathbf{A}})\mathbf{A}_0)}\right)}{\ln(1+\eta\lambda_{r^*}(\mathbf{G}^\natural))} = \frac{\ln\left(\frac{\lambda_1(\mathbf{G}^\natural)}{3\|\mathbf{A}_0\|_{op}^2}\right)}{3\ln(1+\eta\lambda_1(\mathbf{G}^\natural))} = t^* \\ \Leftrightarrow & \ln\left(\frac{8\|\mathbf{A}_0\|_{op}}{\theta\lambda_{\min}(\mathbf{U}_{r^*}^\top(\mathbf{P}_t^{\mathbf{A}})\mathbf{A}_0)}\right) = \frac{\ln(1+\eta\lambda_{r^*}(\mathbf{G}^\natural))}{3\ln(1+\eta\lambda_1(\mathbf{G}^\natural))} \ln\left(\frac{\lambda_1(\mathbf{G}^\natural)}{3\|\mathbf{A}_0\|_{op}^2}\right) \\ \Leftrightarrow & \frac{8\|\mathbf{A}_0\|_{op}}{\theta\lambda_{\min}(\mathbf{U}_{r^*}^\top(\mathbf{P}_t^{\mathbf{A}})\mathbf{A}_0)} = \left(\frac{\lambda_1(\mathbf{G}^\natural)}{3\|\mathbf{A}_0\|_{op}^2}\right)^{\frac{\ln(1+\eta\lambda_{r^*}(\mathbf{G}^\natural))}{3\ln(1+\eta\lambda_1(\mathbf{G}^\natural))}} \\ \Leftrightarrow & \theta = \frac{8\|\mathbf{A}_0\|_{op}}{\lambda_{\min}(\mathbf{U}_{r^*}^\top(\mathbf{P}_t^{\mathbf{A}})\mathbf{A}_0)} \left(\frac{3\|\mathbf{A}_0\|_{op}^2}{\lambda_1(\mathbf{G}^\natural)}\right)^{\frac{\ln(1+\eta\lambda_{r^*}(\mathbf{G}^\natural))}{3\ln(1+\eta\lambda_1(\mathbf{G}^\natural))}} \\ & = \frac{8\|\mathbf{T}\|_{op}}{\lambda_{\min}(\mathbf{U}_{r^*}^\top(\mathbf{P}_t^{\mathbf{A}})\mathbf{T})} \left(\frac{3\|\mathbf{A}_0\|_{op}^2}{\lambda_1(\mathbf{G}^\natural)}\right)^\iota \quad \left[\text{by setting } \iota := \frac{\ln(1+\eta\lambda_{r^*}(\mathbf{G}^\natural))}{3\ln(1+\eta\lambda_1(\mathbf{G}^\natural))} \right] \\ & = \frac{8\|\mathbf{T}\|_{op}}{\lambda_{\min}(\mathbf{U}_{r^*}^\top(\mathbf{P}_t^{\mathbf{A}})\mathbf{T})} \left(\frac{3\|\mathbf{T}\|_{op}^2}{\lambda_1(\mathbf{G}^\natural)}\right)^\iota \alpha^{2\iota}. \end{aligned}$$

In the next, we will discuss how to pick up α . According to Lemma E.3, we need to consider the following two cases on the relationship between r^* and r .

Case 1. $r^* \leq r < 2r^*$: by Lemma E.2 and Lemma E.3, with probability at least $1 - C_1 \exp(-d) - (C_2\xi)^{r-r^*+1} - C_3 \exp(-r)$ for some positive constants C_1, C_2, C_3 , we have

$$\frac{\|\mathbf{T}\|_{op}}{3\sqrt{d}} \leq 1, \quad \frac{\xi}{r\lambda_{\min}(\mathbf{U}_{r^*}^\top(\mathbf{P}_t^{\mathbf{A}})\mathbf{T})} \lesssim 1. \quad (27)$$

Here we pick

$$\alpha \leq \left(\frac{\theta\xi}{24r\sqrt{d}}\right)^{\frac{3\kappa^\natural}{2}} \sqrt{\frac{\lambda_1(\mathbf{G}^\natural)}{27d}},$$

then recall Lemma C.8 on the alignment, we take α here

$$\begin{aligned}
 & \left\| \mathbf{U}_{r^*, \perp}^\top \left(-\nabla_{\mathbf{W}} \tilde{L}(\mathbf{W}^\natural) \right) \mathbf{U}_{r^*}(\mathbf{A}_{t^*}) \right\|_{op} \\
 & \leq \frac{8\|\mathbf{T}\|_{op}}{\lambda_{\min}(\mathbf{U}_{r^*}^\top(\mathbf{P}_t^{\mathbf{A}})\mathbf{T})} \left(\frac{3\|\mathbf{T}\|_{op}^2}{\lambda_1(\mathbf{G}^\natural)} \right)^\iota \alpha^{2\iota} \\
 & = \frac{8\|\mathbf{T}\|_{op}}{\lambda_{\min}(\mathbf{U}_{r^*}^\top(\mathbf{P}_t^{\mathbf{A}})\mathbf{T})} \left(\frac{3\|\mathbf{T}\|_{op}^2}{\lambda_1(\mathbf{G}^\natural)} \right)^\iota \left(\frac{\theta\xi}{24r\sqrt{d}} \right)^{3\kappa^\natural\iota} \left(\frac{\lambda_1(\mathbf{G}^\natural)}{27d} \right)^\iota \\
 & = \frac{8\|\mathbf{T}\|_{op}}{\lambda_{\min}(\mathbf{U}_{r^*}^\top(\mathbf{P}_t^{\mathbf{A}})\mathbf{T})} \left(\frac{\|\mathbf{T}\|_{op}^2}{9d} \right)^\iota \left(\frac{\theta\xi}{24r\sqrt{d}} \right)^{3\kappa^\natural\iota} \\
 & \leq \frac{\|\mathbf{T}\|_{op}\theta\xi}{3r\sqrt{d}\lambda_{\min}(\mathbf{U}_{r^*}^\top(\mathbf{P}_t^{\mathbf{A}})\mathbf{T})} \left(\frac{\|\mathbf{T}\|_{op}^2}{9d} \right)^\iota. \quad \left[\text{since } \iota \geq 1/3\kappa^\natural \text{ and } \frac{\theta\xi}{24r\sqrt{d}} \in (0, 1) \right]
 \end{aligned}$$

Then using Eq. (27), with probability at least $1 - C_1 \exp(-d) - (C_2\xi)^{r-r^*+1} - C_3 \exp(-r)$ for some positive constants C_1, C_2, C_3 , we have

$$\left\| \mathbf{U}_{r^*, \perp}^\top \left(-\nabla_{\mathbf{W}} \tilde{L}(\mathbf{W}^\natural) \right) \mathbf{U}_{r^*}(\mathbf{A}_{t^*}) \right\|_{op} \lesssim \theta.$$

And we can compute the upper bound of t^* as

$$t^* = \frac{\ln \left(\frac{8\|\mathbf{A}\|_{op}}{\theta\lambda_{\min}(\mathbf{U}_{r^*}^\top(\mathbf{P}_t^{\mathbf{A}})\mathbf{A})} \right)}{\ln(1 + \eta\lambda_{r^*}(\mathbf{G}^\natural))} \lesssim \frac{\ln \left(\frac{24r\sqrt{d}}{\theta\xi} \right)}{\ln(1 + \eta\lambda_{r^*}(\mathbf{G}^\natural))}.$$

Case 2. $r \geq 2r^*$: by Lemma E.2 and Lemma E.3, with probability at least $1 - C_4 \exp(-d) - C_5 \exp(-r)$ for some positive constants C_4, C_5 , we have

$$\frac{\|\mathbf{T}\|_{op}}{3\sqrt{d}} \leq 1, \quad \frac{1}{\lambda_{\min}(\mathbf{U}_{r^*}^\top(\mathbf{P}_t^{\mathbf{A}})\mathbf{T})} \lesssim 1.$$

Here we pick

$$\alpha \leq \left(\frac{\theta}{24\sqrt{d}} \right)^{\frac{3\kappa^\natural}{2}} \sqrt{\frac{\lambda_1(\mathbf{G}^\natural)}{27d}}.$$

Similarly, we can obtain

$$\left\| \mathbf{U}_{r^*, \perp}^\top \left(-\nabla_{\mathbf{W}} \tilde{L}(\mathbf{W}^\natural) \right) \mathbf{U}_{r^*}(\mathbf{A}_t) \right\|_{op} \leq \frac{\|\mathbf{T}\|_{op}\theta}{3\sqrt{d}\lambda_{\min}(\mathbf{U}_{r^*}^\top(\mathbf{P}_t^{\mathbf{A}})\mathbf{T})} \left(\frac{\|\mathbf{T}\|_{op}^2}{9d} \right)^\iota \lesssim \theta.$$

And we can compute the upper bound of t^* as

$$t^* \leq \frac{\ln \left(\frac{24\sqrt{d}}{\theta} \right)}{\ln(1 + \eta\lambda_{r^*}(\mathbf{G}^\natural))}.$$

□

Theorem C.10. Under assumptions in Section 2.1 for the linear setting, using the LoRA initialization for $\mathbf{B}_0 = \mathbf{0}$, then for any time-step $t \in \mathbb{N}_+$, we have

$$\left\| \mathbf{V}_{r^*, \perp}^\top \left(-\nabla_{\mathbf{W}} \tilde{L}(\mathbf{W}^\natural) \right) \mathbf{V}_{r^*}(\mathbf{B}_t) \right\|_{op} = 0.$$

Proof. We prove by induction. Recall the complete SVD of Δ in Eq. (1) as

$$\Delta = \tilde{U} \tilde{S}^* \tilde{V}^\top = [U \quad U_\perp] \begin{bmatrix} S^* & \mathbf{0}_{r^* \times (d-r^*)} \\ \mathbf{0}_{(d-r^*) \times r^*} & \mathbf{0}_{(d-r^*) \times (d-r^*)} \end{bmatrix} \begin{bmatrix} V^\top \\ V_\perp^\top \end{bmatrix}.$$

For $t = 1$, recall $G^\natural = \frac{1}{N} \tilde{X}^\top \tilde{X} \Delta$ in Eq. (5), we have

$$B_1 V_\perp = \frac{\eta}{N} A_0^\top G^\natural V_\perp = \frac{\eta}{N} A_0^\top \tilde{X}^\top \tilde{X} \Delta V_\perp = \mathbf{0}_{r \times (d-r^*)}.$$

Assume $B_t V_\perp = \mathbf{0}_{r \times (d-r^*)}$ holds for any $t \in \mathbb{N}_+$ and $t \geq 2$, then

$$B_{t+1} V_\perp = B_t V_\perp - \frac{\eta}{N} A_t^\top \tilde{X}^\top \tilde{X} A_t B_t V_\perp + \frac{\eta}{N} A_t^\top G^\natural V_\perp = \mathbf{0}_{r \times (d-r^*)},$$

which completes the claim. \square

C.2. Gradient Descent under Spectral Initialization

For notational simplicity, we denote $\hat{\Sigma} := \frac{1}{N} \tilde{X}^\top \tilde{X}$ in the following content. Recall the negative gradient of Full Fine-tuning at the first step in Eq. (5), we write it here again

$$G^\natural = -\nabla_{\mathbf{W}} \tilde{L}(\mathbf{W}^\natural) = \frac{1}{N} \tilde{X}^\top \tilde{Y}_\Delta = \hat{\Sigma} \Delta = \tilde{U}_{G^\natural} \tilde{S}_{G^\natural} \tilde{V}_{G^\natural}^\top. \quad (28)$$

In this section, according to Lemma E.1, the following statement

$$\left\| \hat{\Sigma} - I_d \right\|_{op} = \epsilon \leq \min \left\{ \frac{1}{2\kappa}, \frac{c}{\kappa^3} \right\} \leq \frac{1}{2}, \quad \text{for some small constant } c, \quad (29)$$

holds with probability at least $1 - 2C \exp(-\epsilon^2 N)$ for a universal constant $C > 0$. We propose the following initialization scheme (**Spectral-init**)

$$A_0 = \left[\tilde{U}_{G^\natural} \right]_{[:,1:r]} \left[\tilde{S}_{G^\natural}^{1/2} \right]_{[1:r]}, \quad B_0 = \left[\tilde{S}_{G^\natural}^{1/2} \right]_{[1:r]} \left[\tilde{V}_{G^\natural} \right]_{[:,1:r]}^\top.$$

First, we have the following lemma.

Lemma C.11. *Under assumptions in Section 2.1 for the linear setting, with spectral initialization (**Spectral-init**), recall $\kappa := \lambda_1^*(\Delta)/\lambda_{r^*}^*(\Delta)$, then with probability at least with probability $1 - 2C \exp(-\epsilon^2 N)$ for a universal constant $C > 0$, we have*

$$\|A_0 B_0 - \Delta\|_{op} \leq \epsilon \|\Delta\|_{op} \leq \frac{\lambda_{r^*}^*}{2}, \quad (30)$$

and

$$\lambda_{r^*}(A_0) \geq \frac{\sqrt{\lambda_{r^*}^*}}{2}, \quad \lambda_{r^*}(B_0) \geq \frac{\sqrt{\lambda_{r^*}^*}}{2}. \quad (31)$$

Proof. Due to $\text{rank}(G^\natural) = r^*$ and $r \geq r^*$, then $A_0 B_0 = G^\natural$. Accordingly, by Eq. (29), with probability at least $1 - 2 \exp(-c\epsilon^2 N)$, we have

$$\begin{aligned} \|A_0 B_0 - \Delta\|_{op} &\leq \|A_0 B_0 - G^\natural\|_{op} + \|G^\natural - \Delta\|_{op} \\ &= \|G^\natural - \Delta\|_{op} \\ &= \left\| \left(\hat{\Sigma} - I_d \right) \Delta \right\|_{op} && \text{[using Eq. (28)]} \\ &\leq \left\| \hat{\Sigma} - I_d \right\|_{op} \|\Delta\|_{op} \\ &\leq \epsilon \|\Delta\|_{op} \\ &\leq \frac{1}{2\kappa} \|\Delta\|_{op} && \text{[using Eq. (29)]} \\ &= \frac{\lambda_{r^*}^*}{2}. \end{aligned}$$

Then, using the above result and Weyl's inequality, we have the upper bound $\lambda_{r^*}(\mathbf{A}_0\mathbf{B}_0) \leq \lambda_1(\mathbf{A}_0)\lambda_{r^*}(\mathbf{B}_0)$ and the lower bound

$$\lambda_{r^*}(\mathbf{A}_0\mathbf{B}_0) = \lambda_{r^*}(\mathbf{G}^\natural) \geq \lambda_{r^*}(\Delta) - \|\mathbf{G}^\natural - \Delta\|_{op} = \lambda_{r^*}(\Delta) - \|\mathbf{A}_0\mathbf{B}_0 - \Delta\|_{op} \geq \frac{\lambda_{r^*}^*}{2}.$$

Now we are ready to give the lower bound of $\lambda_{r^*}(\mathbf{B}_0)$. Because of $\mathbf{A}_0\mathbf{B}_0 = \mathbf{G}^\natural$ under spectral initialization, we have

$$\lambda_1(\mathbf{A}_0) \leq \sqrt{\lambda_1(\mathbf{G}^\natural)} \leq \sqrt{\|\widehat{\Sigma} - \mathbf{I}_d\|_{op}} \lambda_1(\Delta) \leq \sqrt{\epsilon\lambda_1(\Delta)}, \quad \text{with high probability at least } 1 - 2C \exp(-\epsilon^2 N).$$

where we use $\mathbf{G}^\natural = \widehat{\Sigma}\Delta$ and the concentration results on $\widehat{\Sigma}$. Then combining the above two inequalities, $\lambda_{r^*}(\mathbf{B}_0)$ is lower bounded by

$$\lambda_{r^*}(\mathbf{B}_0) \geq \frac{\lambda_{r^*}(\mathbf{A}_0\mathbf{B}_0)}{\lambda_1(\mathbf{A}_0)} \geq \frac{\lambda_{r^*}^*/2}{\lambda_1(\mathbf{A}_0)} \geq \frac{\sqrt{\lambda_{r^*}^*}}{2},$$

by taking $\epsilon \leq \frac{1}{2\kappa}$. The lower bound of $\lambda_{r^*}(\mathbf{A}_0)$ can be obtained similarly. \square

The following lemma indicates \mathbf{B}_t 's GD dynamics stay in the low-dimensional target subspace under the spectral initialization.

Lemma C.12. *Under assumptions in Section 2.1 for the linear setting, with spectral initialization (Spectral-init), during the iteration, for any $t \in \mathbb{N}^+$, we always have $\mathbf{B}_t\mathbf{V}_\perp = \mathbf{0}_{d \times (d-r^*)}$, where \mathbf{V}_\perp comes from the complete SVD of Δ in Eq. (1).*

Proof. We prove it by induction. First, recall the SVD of Δ in Eq. (1), we have

$$\mathbf{G}^\natural\mathbf{V}_\perp = \widetilde{\Sigma}\Delta\mathbf{V}_\perp = \mathbf{0}_{d \times (d-r^*)},$$

and

$$\begin{aligned} \mathbf{B}_0\mathbf{V}_\perp &= \begin{bmatrix} \widetilde{\mathbf{S}}^{1/2} \\ \mathbf{G}^\natural \end{bmatrix}_{[1:r]} \begin{bmatrix} \widetilde{\mathbf{V}}^\top \\ \mathbf{G}^\natural \end{bmatrix}_{[:,1:r]} \mathbf{V}_\perp \\ &= \begin{bmatrix} \widetilde{\mathbf{S}}^{-1/2} \\ \mathbf{G}^\natural \end{bmatrix}_{[1:r]} \begin{bmatrix} \widetilde{\mathbf{U}}^\top \\ \mathbf{G}^\natural \end{bmatrix}_{[:,1:r]} \mathbf{G}^\natural\mathbf{V}_\perp \\ &= \begin{bmatrix} \widetilde{\mathbf{S}}^{-1/2} \\ \mathbf{G}^\natural \end{bmatrix}_{[1:r]} \begin{bmatrix} \widetilde{\mathbf{U}}^\top \\ \mathbf{G}^\natural \end{bmatrix}_{[:,1:r]} \widehat{\Sigma}\Delta\mathbf{V}_\perp \\ &= \mathbf{0}_{d \times (d-r^*)}. \end{aligned}$$

Next, We prove by induction. Starting from $t = 1$, using the above two equations, we have

$$\begin{aligned} \mathbf{B}_1\mathbf{V}_\perp &= \mathbf{B}_0\mathbf{V}_\perp - \frac{\eta_2}{N} \mathbf{A}_0^\top \widetilde{\mathbf{X}}^\top \left(\widetilde{\mathbf{X}}(\mathbf{W}^\natural + \mathbf{A}_0\mathbf{B}_0) - \widetilde{\mathbf{Y}} \right) \mathbf{V}_\perp \\ &= \mathbf{B}_0\mathbf{V}_\perp - \frac{\eta_2}{N} \mathbf{A}_0^\top \widetilde{\mathbf{X}}^\top \widetilde{\mathbf{X}} \mathbf{A}_0\mathbf{B}_0\mathbf{V}_\perp + \eta \mathbf{A}_0^\top \mathbf{G}^\natural\mathbf{V}_\perp \\ &= \mathbf{0}_{d \times (d-r^*)}. \end{aligned}$$

Assume $\mathbf{B}_t\mathbf{V}_\perp = \mathbf{0}_{d \times (d-r^*)}$ holds for any $t = 2, 3, \dots$, then at $t + 1$, we have

$$\begin{aligned} \mathbf{B}_{t+1}\mathbf{V}_\perp &= \mathbf{B}_t\mathbf{V}_\perp - \frac{\eta}{N} \mathbf{A}_t^\top \widetilde{\mathbf{X}}^\top \widetilde{\mathbf{X}} \mathbf{A}_t\mathbf{B}_t\mathbf{V}_\perp + \eta_2 \mathbf{A}_t^\top \mathbf{G}^\natural\mathbf{V}_\perp \\ &= \mathbf{0}_{d \times (d-r^*)}. \end{aligned}$$

Accordingly we finish the proof. \square

Under spectral initialization, we have already demonstrated that $\mathbf{A}_0\mathbf{B}_0$ is close to Δ . In the following content, we aim to track how $\|\mathbf{A}_t\mathbf{B}_t - \Delta\|_{op}$ behaves (in a local sense), which is a critical ingredient to study both the loss and risk of LoRA training. In this regime, there is no significant difference on setting different step-size η_1 and η_2 . For ease of description, we set $\eta_1 = \eta_2 := \eta$.

Here we can characterize the operator norm of $(\mathbf{A}_t \mathbf{B}_t - \Delta)$ as

$$\begin{aligned}
 \|\mathbf{A}_t \mathbf{B}_t - \Delta\|_{op} &= \left\| \left(\mathbf{A}_t \mathbf{B}_t - \Delta \right) \begin{bmatrix} \mathbf{V} & \mathbf{V}_\perp \end{bmatrix} \right\|_{op} && \text{[by unitary invariance of operator norm]} \\
 &= \|\mathbf{A}_t \mathbf{B}_t \mathbf{V} - \mathbf{U} \mathbf{S}^*\|_{op} && \text{[by Lemma C.12]} \\
 &= \left\| \left(\mathbf{U} \mathbf{U}^\top + \mathbf{U}_\perp \mathbf{U}_\perp^\top \right) \left(\mathbf{A}_t \mathbf{B}_t \mathbf{V} - \mathbf{U} \mathbf{S}^* \right) \right\|_{op} \\
 &= \left\| \mathbf{U} \left(\mathbf{U}^\top \mathbf{A}_t \mathbf{B}_t \mathbf{V} - \mathbf{S}^* \right) \right\|_{op} + \|\mathbf{U}_\perp \mathbf{U}_\perp^\top \mathbf{A}_t \mathbf{B}_t \mathbf{V}\|_{op} \\
 &\leq \underbrace{\|\mathbf{U}^\top \mathbf{A}_t \mathbf{B}_t \mathbf{V} - \mathbf{S}^*\|_{op}}_{\text{signal space}} + \underbrace{\|\mathbf{U}_\perp^\top \mathbf{A}_t \mathbf{B}_t \mathbf{V}\|_{op}}_{\text{complementary}}, \tag{32}
 \end{aligned}$$

where the first term denotes the loss in the signal space $\|\mathbf{U}^\top \mathbf{A} \mathbf{B} \mathbf{V} - \mathbf{S}^*\|_{op}$ and the second term denotes the complementary space decay $\|\mathbf{U}_\perp^\top \mathbf{A} \mathbf{B} \mathbf{V}\|_{op}$. Next, we need a new parametrization to track the dynamics of these two terms. Recall the complete SVD of Δ in Eq. (1) as

$$\Delta = \tilde{\mathbf{U}} \tilde{\mathbf{S}} \tilde{\mathbf{V}}^\top = \begin{bmatrix} \mathbf{U} & \mathbf{U}_\perp \end{bmatrix} \begin{bmatrix} \mathbf{S}^* & \mathbf{0}_{r^* \times (d-r^*)} \\ \mathbf{0}_{(d-r^*) \times r^*} & \mathbf{0}_{(d-r^*) \times (d-r^*)} \end{bmatrix} \begin{bmatrix} \mathbf{V}^\top \\ \mathbf{V}_\perp^\top \end{bmatrix}.$$

For notational simplicity, we denote

$$\mathbf{A}_t^U := \mathbf{U}^\top \mathbf{A}_t, \quad \mathbf{A}_t^{U_\perp} := \mathbf{U}_\perp^\top \mathbf{A}_t, \quad \mathbf{B}_t \mathbf{V} := \mathbf{B}_t^V, \quad \mathbf{B}_t \mathbf{V}_\perp := \mathbf{B}_t^{V_\perp}.$$

and thus

$$\mathbf{R}_t := (\mathbf{A}_t \mathbf{B}_t - \Delta) \mathbf{V}, \quad \mathbf{R}_t^* := \mathbf{A}_t^U \mathbf{B}_t^V - \mathbf{S}^*, \quad \mathbf{R}_t^\perp := \mathbf{A}_t^{U_\perp} \mathbf{B}_t^{V_\perp}.$$

Accordingly, Eq. (32) can be reformulated as $\|\mathbf{R}_t\|_{op} \leq \|\mathbf{R}_t^*\|_{op} + \|\mathbf{R}_t^\perp\|_{op}$. By Lemma C.12, we have $\mathbf{B}^{V_\perp} = \mathbf{0}_{r \times (k-r^*)}$ $\forall t \in \mathbb{N}^+$. Next, we can track \mathbf{R}_t^* and \mathbf{R}_t^\perp via the following two lemmas.

Lemma C.13. *Under assumptions in Section 2.1 for the linear setting, with spectral initialization (Spectral-init), we have the following reparametrized iterates*

$$\mathbf{A}_{t+1}^U = \mathbf{A}_t^U - \eta \mathbf{R}_t^* (\mathbf{B}_t^V)^\top - \eta \mathbf{U}^\top (\hat{\Sigma} - \mathbf{I}_d) \mathbf{R}_t (\mathbf{B}_t^V)^\top, \tag{33}$$

$$\mathbf{A}_{t+1}^{U_\perp} = \mathbf{A}_t^{U_\perp} - \eta \mathbf{R}_t^\perp (\mathbf{B}_t^{V_\perp})^\top - \eta \mathbf{U}_\perp^\top (\hat{\Sigma} - \mathbf{I}_d) \mathbf{R}_t (\mathbf{B}_t^{V_\perp})^\top, \tag{34}$$

$$\begin{aligned}
 \mathbf{B}_{t+1}^V &= \mathbf{B}_t^V - \eta (\mathbf{A}_t^U)^\top \mathbf{R}_t^* - \eta (\mathbf{A}_t^{U_\perp})^\top \mathbf{R}_t^\perp \\
 &\quad - \eta (\mathbf{A}_t^U)^\top \mathbf{U}^\top (\hat{\Sigma} - \mathbf{I}_d) \mathbf{R}_t - \eta (\mathbf{A}_t^{U_\perp})^\top \mathbf{U}_\perp^\top (\hat{\Sigma} - \mathbf{I}_d) \mathbf{R}_t. \tag{35}
 \end{aligned}$$

Proof. Recall the gradient update for \mathbf{A}_{t+1} , we have

$$\begin{aligned}
 \mathbf{A}_{t+1} &= \mathbf{A}_t - \eta \hat{\Sigma} (\mathbf{A}_t \mathbf{B}_t - \Delta) (\mathbf{B}_t)^\top \\
 &= \mathbf{A}_t - \eta (\mathbf{A}_t \mathbf{B}_t - \Delta) (\mathbf{B}_t)^\top - \eta (\hat{\Sigma} - \mathbf{I}_d) (\mathbf{A}_t \mathbf{B}_t - \Delta) (\mathbf{B}_t)^\top.
 \end{aligned}$$

Recall $\mathbf{R}_t := (\mathbf{A}_t \mathbf{B}_t - \Delta) \mathbf{V}$ and $\Delta = \mathbf{U} \mathbf{S}^* \mathbf{V}^\top$, we have

$$\begin{aligned}
 \mathbf{U}^\top \mathbf{A}_{t+1} &= \mathbf{U}^\top \mathbf{A}_t - \eta \mathbf{U}^\top (\mathbf{A}_t \mathbf{B}_t - \Delta) (\mathbf{V} \mathbf{V}^\top + \mathbf{V}_\perp \mathbf{V}_\perp^\top) (\mathbf{B}_t)^\top \\
 &\quad - \eta \mathbf{U}^\top (\hat{\Sigma} - \mathbf{I}_d) (\mathbf{A}_t \mathbf{B}_t - \Delta) (\mathbf{V} \mathbf{V}^\top + \mathbf{V}_\perp \mathbf{V}_\perp^\top) (\mathbf{B}_t)^\top \\
 &= \mathbf{U}^\top \mathbf{A}_t - \eta \mathbf{U}^\top (\mathbf{A}_t \mathbf{B}_t \mathbf{V} - \Delta \mathbf{V}) (\mathbf{B}_t \mathbf{V})^\top - \eta \mathbf{U}^\top (\hat{\Sigma} - \mathbf{I}_d) (\mathbf{A}_t \mathbf{B}_t \mathbf{V} - \Delta \mathbf{V}) (\mathbf{B}_t \mathbf{V})^\top \\
 &\quad \text{[by Lemma C.12]} \\
 &= \mathbf{U}^\top \mathbf{A}_t - \eta (\mathbf{U}^\top \mathbf{A}_t \mathbf{B}_t \mathbf{V} - \mathbf{S}^*) (\mathbf{B}_t \mathbf{V})^\top - \eta \mathbf{U}^\top (\hat{\Sigma} - \mathbf{I}_d) \mathbf{R}_t (\mathbf{B}_t \mathbf{V})^\top.
 \end{aligned}$$

Accordingly, the recursion for \mathbf{A}_{t+1}^U is reformulated as

$$\mathbf{A}_{t+1}^U = \mathbf{A}_t^U - \eta \mathbf{R}_t^* (\mathbf{B}_t^V)^\top - \eta \mathbf{U}^\top (\hat{\Sigma} - \mathbf{I}_d) \mathbf{R}_t (\mathbf{B}_t^V)^\top.$$

Similarly, we can obtain

$$\mathbf{A}_{t+1}^{U^\perp} = \mathbf{A}_t^{U^\perp} - \eta \mathbf{R}_t^\perp (\mathbf{B}_t^V)^\top - \eta \mathbf{U}_\perp^\top (\hat{\Sigma} - \mathbf{I}_d) \mathbf{R}_t (\mathbf{B}_t^V)^\top.$$

Regarding the recursion for \mathbf{B}_{t+1} , we can derive in a similar way

$$\begin{aligned} \mathbf{B}_{t+1} \mathbf{V} &= \mathbf{B}_t \mathbf{V} - \eta (\mathbf{A}_t)^\top \hat{\Sigma} (\mathbf{A}_t \mathbf{B}_t - \Delta) \mathbf{V} \\ &= \mathbf{B}_t \mathbf{V} - \eta (\mathbf{A}_t)^\top (\mathbf{U} \mathbf{U}^\top + \mathbf{U}_\perp \mathbf{U}_\perp^\top) (\mathbf{A}_t \mathbf{B}_t - \Delta) \mathbf{V} \\ &\quad - \eta (\mathbf{A}_t)^\top (\mathbf{U} \mathbf{U}^\top + \mathbf{U}_\perp \mathbf{U}_\perp^\top) (\hat{\Sigma} - \mathbf{I}_d) (\mathbf{A}_t \mathbf{B}_t - \Delta) \mathbf{V}, \end{aligned}$$

which implies

$$\mathbf{B}_{t+1}^V = \mathbf{B}_t^V - \eta (\mathbf{A}_t^U)^\top \mathbf{R}_t^* - \eta (\mathbf{A}_t^{U^\perp})^\top \mathbf{R}_t^\perp - \eta (\mathbf{A}_t^U)^\top \mathbf{U}^\top (\hat{\Sigma} - \mathbf{I}_d) \mathbf{R}_t - \eta (\mathbf{A}_t^{U^\perp})^\top \mathbf{U}_\perp^\top (\hat{\Sigma} - \mathbf{I}_d) \mathbf{R}_t.$$

□

In the next, we are able to characterize the upper bound of $\|\mathbf{R}_{t+1}^*\|_{op}$.

Lemma C.14. Denote $\mathcal{M}_t := \max \left\{ \|\mathbf{R}_t^*\|_{op}, \|\mathbf{R}_t^\perp\|_{op} \right\}$, under assumptions in Section 2.1 for the linear setting, with spectral initialization (Spectral-init), then we choose ϵ with probability at least $1 - 2C \exp(-\epsilon^2 N)$ for a universal constant $C > 0$, we have

$$\begin{aligned} \|\mathbf{R}_{t+1}^*\|_{op} &\leq \left(1 - \eta (\lambda_{r^*}^2 (\mathbf{A}_t^U) + \lambda_{r^*}^2 (\mathbf{B}_t^V)) \right) \mathcal{M}_t \\ &\quad + 2\eta\epsilon \|\mathbf{B}_t^V\|_{op}^2 \mathcal{M}_t + \eta^2 \|\mathbf{A}_t^U\|_{op} \|\mathbf{B}_t^V\|_{op} \mathcal{M}_t^2 + 2\eta^2\epsilon \|\mathbf{A}_t^U\|_{op} \|\mathbf{B}_t^V\|_{op} \mathcal{M}_t^2 \\ &\quad + \eta \|\mathbf{A}_t^U\|_{op} \|\mathbf{A}_t^{U^\perp}\|_{op} \mathcal{M}_t + \eta^2 \|\mathbf{B}_t^V\|_{op} \|\mathbf{A}_t^{U^\perp}\|_{op} \mathcal{M}_t^2 \\ &\quad + 2\eta^2\epsilon \|\mathbf{B}_t^V\|_{op} \|\mathbf{A}_t^{U^\perp}\|_{op} \mathcal{M}_t^2 + 2\eta\epsilon \|\mathbf{A}_t^U\|_{op}^2 \mathcal{M}_t \\ &\quad + 2\eta^2\epsilon \|\mathbf{A}_t^U\|_{op} \|\mathbf{B}_t^V\|_{op} \mathcal{M}_t + 4\eta^2\epsilon^2 \|\mathbf{A}_t^U\|_{op} \|\mathbf{B}_t^V\|_{op} \mathcal{M}_t^2 \\ &\quad + 2\eta\epsilon \|\mathbf{A}_t^U\|_{op} \|\mathbf{A}_t^{U^\perp}\|_{op} \mathcal{M}_t + 2\eta^2\epsilon \|\mathbf{A}_t^{U^\perp}\|_{op} \|\mathbf{B}_t^V\|_{op} \mathcal{M}_t^2 \\ &\quad + 4\eta^2\epsilon^2 \|\mathbf{A}_t^{U^\perp}\|_{op} \|\mathbf{B}_t^V\|_{op} \mathcal{M}_t^2, \end{aligned}$$

and

$$\begin{aligned} \|\mathbf{R}_{t+1}^\perp\|_{op} &\leq \left(1 - \eta (\lambda_{\min}^2 (\mathbf{A}_t^{U^\perp}) + \lambda_{r^*}^2 (\mathbf{B}_t^V)) \right) \mathcal{M}_t \\ &\quad + 2\eta\epsilon \|\mathbf{B}_t^V\|_{op}^2 \mathcal{M}_t + \eta^2 \|\mathbf{A}_t^U\|_{op} \|\mathbf{B}_t^V\|_{op} \mathcal{M}_t^2 + 2\eta^2\epsilon \|\mathbf{A}_t^U\|_{op} \|\mathbf{B}_t^V\|_{op} \mathcal{M}_t^2 \\ &\quad + \eta \|\mathbf{A}_t^U\|_{op} \|\mathbf{A}_t^{U^\perp}\|_{op} \mathcal{M}_t + \eta^2 \|\mathbf{B}_t^V\|_{op} \|\mathbf{A}_t^{U^\perp}\|_{op} \mathcal{M}_t^2 \\ &\quad + 2\eta^2\epsilon \|\mathbf{B}_t^V\|_{op} \|\mathbf{A}_t^{U^\perp}\|_{op} \mathcal{M}_t^2 + 2\eta\epsilon \|\mathbf{A}_t^U\|_{op} \|\mathbf{A}_t^{U^\perp}\|_{op} \mathcal{M}_t \\ &\quad + 2\eta^2\epsilon \|\mathbf{A}_t^U\|_{op} \|\mathbf{B}_t^V\|_{op} \mathcal{M}_t + 4\eta^2\epsilon^2 \|\mathbf{A}_t^U\|_{op} \|\mathbf{B}_t^V\|_{op} \mathcal{M}_t^2 \\ &\quad + 2\eta\epsilon \|\mathbf{A}_t^{U^\perp}\|_{op}^2 \mathcal{M}_t + 2\eta^2\epsilon \|\mathbf{A}_t^{U^\perp}\|_{op} \|\mathbf{B}_t^V\|_{op} \mathcal{M}_t^2 \\ &\quad + 4\eta^2\epsilon^2 \|\mathbf{A}_t^{U^\perp}\|_{op} \|\mathbf{B}_t^V\|_{op} \mathcal{M}_t^2. \end{aligned} \tag{36}$$

Proof. Here we first track the dynamics of \mathbf{R}_t^* . We have

$$\begin{aligned}
 \mathbf{R}_{t+1}^* &= \mathbf{A}_{t+1}^U \mathbf{B}_{t+1}^V - \mathbf{S}^* \\
 &= \mathbf{R}_t^* - \eta \mathbf{R}_t^* (\mathbf{B}_t^V)^\top \mathbf{B}_t^V - \eta \mathbf{U}^\top (\widehat{\Sigma} - \mathbf{I}_d) \mathbf{R}_t (\mathbf{B}_t^V)^\top \mathbf{B}_t^V \\
 &\quad - \eta \mathbf{A}_t^U (\mathbf{A}_t^U)^\top \mathbf{R}_t^* + \eta^2 \mathbf{R}_t^* (\mathbf{B}_t^V)^\top (\mathbf{A}_t^U)^\top \mathbf{R}_t^* + \eta^2 \mathbf{U}^\top (\widehat{\Sigma} - \mathbf{I}_d) \mathbf{R}_t (\mathbf{B}_t^V)^\top (\mathbf{A}_t^U)^\top \mathbf{R}_t^* \\
 &\quad - \eta \mathbf{A}_t^U (\mathbf{A}_t^{U\perp})^\top \mathbf{R}_t^\perp + \eta^2 \mathbf{R}_t^* (\mathbf{B}_t^V)^\top (\mathbf{A}_t^{U\perp})^\top \mathbf{R}_t^\perp + \eta^2 \mathbf{U}^\top (\widehat{\Sigma} - \mathbf{I}_d) \mathbf{R}_t (\mathbf{B}_t^V)^\top (\mathbf{A}_t^{U\perp})^\top \mathbf{R}_t^\perp \\
 &\quad - \eta \mathbf{A}_t^U (\mathbf{A}_t^U)^\top \mathbf{U}^\top (\widehat{\Sigma} - \mathbf{I}_d) \mathbf{R}_t \\
 &\quad + \eta^2 \mathbf{R}_t^* (\mathbf{B}_t^V)^\top (\mathbf{A}_t^U)^\top \mathbf{U}^\top (\widehat{\Sigma} - \mathbf{I}_d) \mathbf{R}_t \\
 &\quad + \eta^2 \mathbf{U}^\top (\widehat{\Sigma} - \mathbf{I}_d) \mathbf{R}_t (\mathbf{B}_t^V)^\top (\mathbf{A}_t^U)^\top \mathbf{U}^\top (\widehat{\Sigma} - \mathbf{I}_d) \mathbf{R}_t \\
 &\quad - \eta \mathbf{A}_t^U (\mathbf{A}_t^{U\perp})^\top \mathbf{U}_\perp^\top (\widehat{\Sigma} - \mathbf{I}_d) \mathbf{R}_t \\
 &\quad + \eta^2 \mathbf{R}_t^* (\mathbf{B}_t^V)^\top (\mathbf{A}_t^{U\perp})^\top \mathbf{U}_\perp^\top (\widehat{\Sigma} - \mathbf{I}_d) \mathbf{R}_t \\
 &\quad + \eta^2 \mathbf{U}^\top (\widehat{\Sigma} - \mathbf{I}_d) \mathbf{R}_t (\mathbf{B}_t^V)^\top (\mathbf{A}_t^{U\perp})^\top \mathbf{U}_\perp^\top (\widehat{\Sigma} - \mathbf{I}_d) \mathbf{R}_t.
 \end{aligned}$$

Then, we take operator norm over the above equation. Hence, with probability at least $1 - 2C \exp(-\epsilon^2 N)$ for a universal constant $C > 0$, we have

$$\begin{aligned}
 \|\mathbf{R}_{t+1}^*\|_{op} &\leq \left(1 - \eta (\lambda_{r^*}^2(\mathbf{A}_t^U) + \lambda_{r^*}^2(\mathbf{B}_t^V))\right) \|\mathbf{R}_t^*\|_{op} \\
 &\quad + \eta \epsilon \|\mathbf{B}_t^V\|_{op}^2 \|\mathbf{R}_t\|_{op} + \eta^2 \|\mathbf{A}_t^U\|_{op} \|\mathbf{B}_t^V\|_{op} \|\mathbf{R}_t^*\|_{op}^2 + \eta^2 \epsilon \|\mathbf{A}_t^U\|_{op} \|\mathbf{B}_t^V\|_{op} \|\mathbf{R}_t^*\|_{op} \|\mathbf{R}_t\|_{op} \\
 &\quad + \eta \|\mathbf{A}_t^U\|_{op} \|\mathbf{A}_t^{U\perp}\|_{op} \|\mathbf{R}_t^\perp\|_{op} + \eta^2 \|\mathbf{B}_t^V\|_{op} \|\mathbf{A}_t^{U\perp}\|_{op} \|\mathbf{R}_t^\perp\|_{op} \|\mathbf{R}_t^*\|_{op} \\
 &\quad + \eta^2 \epsilon \|\mathbf{B}_t^V\|_{op} \|\mathbf{A}_t^{U\perp}\|_{op} \|\mathbf{R}_t^\perp\|_{op} \|\mathbf{R}_t\|_{op} + \eta \epsilon \|\mathbf{A}_t^U\|_{op}^2 \|\mathbf{R}_t\|_{op} \\
 &\quad + \eta^2 \epsilon \|\mathbf{A}_t^U\|_{op} \|\mathbf{B}_t^V\|_{op} \|\mathbf{R}_t\|_{op} + \eta^2 \epsilon^2 \|\mathbf{A}_t^U\|_{op} \|\mathbf{B}_t^V\|_{op} \|\mathbf{R}_t\|_{op}^2 \\
 &\quad + \eta \epsilon \|\mathbf{A}_t^U\|_{op} \|\mathbf{A}_t^{U\perp}\|_{op} \|\mathbf{R}_t\|_{op} + \eta^2 \epsilon \|\mathbf{A}_t^{U\perp}\|_{op} \|\mathbf{B}_t^V\|_{op} \|\mathbf{R}_t^*\|_{op} \|\mathbf{R}_t\|_{op} \\
 &\quad + \eta^2 \epsilon^2 \|\mathbf{A}_t^{U\perp}\|_{op} \|\mathbf{B}_t^V\|_{op} \|\mathbf{R}_t\|_{op}^2.
 \end{aligned}$$

Next, we take maximum over $\|\mathbf{R}_t^*\|_{op}$ and $\|\mathbf{R}_t^\perp\|_{op}$ on the right hand side above. Recall $\mathcal{M}_t = \max\{\|\mathbf{R}_t^*\|_{op}, \|\mathbf{R}_t^\perp\|_{op}\}$,

using the fact that $\|\mathbf{R}_t\|_{op} \leq 2\mathcal{M}_t$, we have:

$$\begin{aligned}
 \|\mathbf{R}_{t+1}^*\|_{op} &\leq \left(1 - \eta (\lambda_{r^*}^2(\mathbf{A}_t^U) + \lambda_{r^*}^2(\mathbf{B}_t^V))\right) \mathcal{M}_t \\
 &\quad + 2\eta\epsilon \|\mathbf{B}_t^V\|_{op}^2 \mathcal{M}_t + \eta^2 \|\mathbf{A}_t^U\|_{op} \|\mathbf{B}_t^V\|_{op} \mathcal{M}_t^2 + 2\eta^2\epsilon \|\mathbf{A}_t^U\|_{op} \|\mathbf{B}_t^V\|_{op} \mathcal{M}_t^2 \\
 &\quad + \eta \|\mathbf{A}_t^U\|_{op} \|\mathbf{A}_t^{U\perp}\|_{op} \mathcal{M}_t + \eta^2 \|\mathbf{B}_t^V\|_{op} \|\mathbf{A}_t^{U\perp}\|_{op} \mathcal{M}_t^2 \\
 &\quad + 2\eta^2\epsilon \|\mathbf{B}_t^V\|_{op} \|\mathbf{A}_t^{U\perp}\|_{op} \mathcal{M}_t^2 + 2\eta\epsilon \|\mathbf{A}_t^U\|_{op}^2 \mathcal{M}_t \\
 &\quad + 2\eta^2\epsilon \|\mathbf{A}_t^U\|_{op} \|\mathbf{B}_t^V\|_{op} \mathcal{M}_t + 4\eta^2\epsilon^2 \|\mathbf{A}_t^U\|_{op} \|\mathbf{B}_t^V\|_{op} \mathcal{M}_t^2 \\
 &\quad + 2\eta\epsilon \|\mathbf{A}_t^U\|_{op} \|\mathbf{A}_t^{U\perp}\|_{op} \mathcal{M}_t + 2\eta^2\epsilon \|\mathbf{A}_t^{U\perp}\|_{op} \|\mathbf{B}_t^V\|_{op} \mathcal{M}_t^2 \\
 &\quad + 4\eta^2\epsilon^2 \|\mathbf{A}_t^{U\perp}\|_{op} \|\mathbf{B}_t^V\|_{op} \mathcal{M}_t^2.
 \end{aligned}$$

Next, we track the dynamics of \mathbf{R}_t^\perp . We have

$$\begin{aligned}
 \mathbf{R}_{t+1}^\perp &= \mathbf{A}_{t+1}^{U\perp} \mathbf{B}_{t+1}^V \\
 &= \mathbf{R}_t^\perp - \eta \mathbf{R}_t^\perp (\mathbf{B}_t^V)^\top \mathbf{B}_t^V - \eta \mathbf{U}_\perp^\top (\widehat{\Sigma} - \mathbf{I}_d) \mathbf{R}_t (\mathbf{B}_t^V)^\top \mathbf{B}_t^V \\
 &\quad - \eta \mathbf{A}_t^{U\perp} (\mathbf{A}_t^U)^\top \mathbf{R}_t^* + \eta^2 \mathbf{R}_t^\perp (\mathbf{B}_t^V)^\top (\mathbf{A}_t^U)^\top \mathbf{R}_t^* + \eta^2 \mathbf{U}_\perp^\top (\widehat{\Sigma} - \mathbf{I}_d) \mathbf{R}_t (\mathbf{B}_t^V)^\top (\mathbf{A}_t^U)^\top \mathbf{R}_t^* \\
 &\quad - \eta \mathbf{A}_t^{U\perp} (\mathbf{A}_t^{U\perp})^\top \mathbf{R}_t^\perp + \eta^2 \mathbf{R}_t^\perp (\mathbf{B}_t^V)^\top (\mathbf{A}_t^{U\perp})^\top \mathbf{R}_t^\perp + \eta^2 \mathbf{U}_\perp^\top (\widehat{\Sigma} - \mathbf{I}_d) \mathbf{R}_t (\mathbf{B}_t^V)^\top (\mathbf{A}_t^{U\perp})^\top \mathbf{R}_t^\perp \\
 &\quad - \eta \mathbf{A}_t^{U\perp} (\mathbf{A}_t^U)^\top \mathbf{U}_\perp^\top (\widehat{\Sigma} - \mathbf{I}_d) \mathbf{R}_t \\
 &\quad + \eta^2 \mathbf{R}_t^\perp (\mathbf{B}_t^V)^\top (\mathbf{A}_t^U)^\top \mathbf{U}_\perp^\top (\widehat{\Sigma} - \mathbf{I}_d) \mathbf{R}_t \\
 &\quad + \eta^2 \mathbf{U}_\perp^\top (\widehat{\Sigma} - \mathbf{I}_d) \mathbf{R}_t (\mathbf{B}_t^V)^\top (\mathbf{A}_t^U)^\top \mathbf{U}_\perp^\top (\widehat{\Sigma} - \mathbf{I}_d) \mathbf{R}_t \\
 &\quad - \eta \mathbf{A}_t^{U\perp} (\mathbf{A}_t^{U\perp})^\top \mathbf{U}_\perp^\top (\widehat{\Sigma} - \mathbf{I}_d) \mathbf{R}_t \\
 &\quad + \eta^2 \mathbf{R}_t^\perp (\mathbf{B}_t^V)^\top (\mathbf{A}_t^{U\perp})^\top \mathbf{U}_\perp^\top (\widehat{\Sigma} - \mathbf{I}_d) \mathbf{R}_t \\
 &\quad + \eta^2 \mathbf{U}_\perp^\top (\widehat{\Sigma} - \mathbf{I}_d) \mathbf{R}_t (\mathbf{B}_t^V)^\top (\mathbf{A}_t^{U\perp})^\top \mathbf{U}_\perp^\top (\widehat{\Sigma} - \mathbf{I}_d) \mathbf{R}_t.
 \end{aligned}$$

Then, we take operator norm over the above equation. With probability at least $1 - 2C \exp(-\epsilon^2 N)$ for a universal constant $C > 0$, we have

$$\begin{aligned}
 \|\mathbf{R}_{t+1}^\perp\|_{op} &\leq \left(1 - \eta (\lambda_{\min}^2(\mathbf{A}_t^{U\perp}) + \lambda_{r^*}^2(\mathbf{B}_t^V))\right) \|\mathbf{R}_t^\perp\|_{op} \\
 &\quad + \eta\epsilon \|\mathbf{B}_t^V\|_{op}^2 \|\mathbf{R}_t\|_{op} + \eta^2 \|\mathbf{A}_t^U\|_{op} \|\mathbf{B}_t^V\|_{op} \|\mathbf{R}_t^*\|_{op} \|\mathbf{R}_t^\perp\|_{op} + \eta^2\epsilon \|\mathbf{A}_t^U\|_{op} \|\mathbf{B}_t^V\|_{op} \|\mathbf{R}_t^*\|_{op} \|\mathbf{R}_t\|_{op} \\
 &\quad + \eta \|\mathbf{A}_t^U\|_{op} \|\mathbf{A}_t^{U\perp}\|_{op} \|\mathbf{R}_t^*\|_{op} + \eta^2 \|\mathbf{B}_t^V\|_{op} \|\mathbf{A}_t^{U\perp}\|_{op} \|\mathbf{R}_t^\perp\|_{op}^2 \\
 &\quad + \eta^2\epsilon \|\mathbf{B}_t^V\|_{op} \|\mathbf{A}_t^{U\perp}\|_{op} \|\mathbf{R}_t^\perp\|_{op} \|\mathbf{R}_t\|_{op} + \eta\epsilon \|\mathbf{A}_t^U\|_{op} \|\mathbf{A}_t^{U\perp}\|_{op} \|\mathbf{R}_t\|_{op} \\
 &\quad + \eta^2\epsilon \|\mathbf{A}_t^U\|_{op} \|\mathbf{B}_t^V\|_{op} \|\mathbf{R}_t\|_{op} + \eta^2\epsilon^2 \|\mathbf{A}_t^U\|_{op} \|\mathbf{B}_t^V\|_{op} \|\mathbf{R}_t\|_{op}^2 \\
 &\quad + \eta\epsilon \|\mathbf{A}_t^{U\perp}\|_{op}^2 \|\mathbf{R}_t\|_{op} + \eta^2\epsilon \|\mathbf{A}_t^{U\perp}\|_{op} \|\mathbf{B}_t^V\|_{op} \|\mathbf{R}_t^\perp\|_{op} \|\mathbf{R}_t\|_{op} \\
 &\quad + \eta^2\epsilon^2 \|\mathbf{A}_t^{U\perp}\|_{op} \|\mathbf{B}_t^V\|_{op} \|\mathbf{R}_t\|_{op}^2.
 \end{aligned}$$

Next, we take maximum over $\|\mathbf{R}_t^*\|_{op}$ and $\|\mathbf{R}_t^\perp\|_{op}$ on the right hand side above. Recall $\mathcal{M}_t = \max\{\|\mathbf{R}_t^*\|_{op}, \|\mathbf{R}_t^\perp\|_{op}\}$, using the fact that $\|\mathbf{R}_t\|_{op} \leq 2\mathcal{M}_t$, we have:

$$\begin{aligned} \|\mathbf{R}_{t+1}^\perp\|_{op} &\leq \left(1 - \eta \left(\lambda_{\min}^2(\mathbf{A}_t^{U^\perp}) + \lambda_{r^*}^2(\mathbf{B}_t^V)\right)\right) \mathcal{M}_t \\ &\quad + 2\eta\epsilon \|\mathbf{B}_t^V\|_{op}^2 \mathcal{M}_t + \eta^2 \|\mathbf{A}_t^U\|_{op} \|\mathbf{B}_t^V\|_{op} \mathcal{M}_t^2 + 2\eta^2\epsilon \|\mathbf{A}_t^U\|_{op} \|\mathbf{B}_t^V\|_{op} \mathcal{M}_t^2 \\ &\quad + \eta \|\mathbf{A}_t^U\|_{op} \|\mathbf{A}_t^{U^\perp}\|_{op} \mathcal{M}_t + \eta^2 \|\mathbf{B}_t^V\|_{op} \|\mathbf{A}_t^{U^\perp}\|_{op} \mathcal{M}_t^2 \\ &\quad + 2\eta^2\epsilon \|\mathbf{B}_t^V\|_{op} \|\mathbf{A}_t^{U^\perp}\|_{op} \mathcal{M}_t^2 + 2\eta\epsilon \|\mathbf{A}_t^U\|_{op} \|\mathbf{A}_t^{U^\perp}\|_{op} \mathcal{M}_t \\ &\quad + 2\eta^2\epsilon \|\mathbf{A}_t^U\|_{op} \|\mathbf{B}_t^V\|_{op} \mathcal{M}_t + 4\eta^2\epsilon^2 \|\mathbf{A}_t^U\|_{op} \|\mathbf{B}_t^V\|_{op} \mathcal{M}_t^2 \\ &\quad + 2\eta\epsilon \|\mathbf{A}_t^{U^\perp}\|_{op}^2 \mathcal{M}_t + 2\eta^2\epsilon \|\mathbf{A}_t^{U^\perp}\|_{op} \|\mathbf{B}_t^V\|_{op} \mathcal{M}_t^2 \\ &\quad + 4\eta^2\epsilon^2 \|\mathbf{A}_t^{U^\perp}\|_{op} \|\mathbf{B}_t^V\|_{op} \mathcal{M}_t^2. \end{aligned}$$

Finally we conclude the proof. \square

Before we move to the main proof, we need to establish a strict upper bound on \mathbf{A}_t and \mathbf{B}_t .

Lemma C.15. *Under assumptions in Section 2.1 for the linear setting, suppose $\|\mathbf{A}_t^\top \mathbf{A}_t - \mathbf{B}_t^\top \mathbf{B}_t\|_{op} + \epsilon \|\mathbf{R}_t\|_{op} \leq \lambda_1^*$ and $\eta \leq \frac{1}{10\lambda_1^*}$, if $\|\mathbf{A}_t\|_{op} \leq 2\sqrt{\lambda_1^*}$ and $\|\mathbf{B}_t\|_{op} \leq 2\sqrt{\lambda_1^*}$, we choose ϵ satisfying Eq. (29), then with probability $1 - 2C \exp(-\epsilon^2 N)$ for a universal constant $C > 0$, we have*

$$\|\mathbf{A}_{t+1}\|_{op} \leq 2\sqrt{\lambda_1^*}, \quad \|\mathbf{B}_{t+1}\|_{op} \leq 2\sqrt{\lambda_1^*}.$$

Proof. Inspired by Soltanolkotabi et al. (2023), we recall the stacked iterate \mathbf{Z}_t defined in Eq. (12) and construct an anti-iterate

$$\underline{\mathbf{Z}}_t := \begin{bmatrix} \mathbf{A}_t \\ -\mathbf{B}_t^\top \end{bmatrix}.$$

Additionally, we define a perturbation matrix

$$\underline{\mathbf{\Xi}}_t := \begin{bmatrix} \mathbf{0}_{d \times d} & (\tilde{\Sigma} - \mathbf{I}_d) \mathbf{R}_t \\ \mathbf{R}_t^\top (\tilde{\Sigma} - \mathbf{I}_d) & \mathbf{0}_{k \times k} \end{bmatrix}.$$

Then, we can reformulate the recursion of \mathbf{Z}_{t+1} as

$$\begin{aligned} \mathbf{Z}_{t+1} &= \mathbf{Z}_t - \eta \left(\mathbf{Z}_t \mathbf{Z}_t^\top - \underline{\mathbf{Z}}_t \underline{\mathbf{Z}}_t^\top - \mathbf{\Gamma} \right) \mathbf{Z}_t + \eta \underline{\mathbf{\Xi}}_t \mathbf{Z}_t \\ &= (\mathbf{I}_{2d} - \eta \mathbf{Z}_t \mathbf{Z}_t^\top) \mathbf{Z}_t + \eta \underline{\mathbf{Z}}_t \underline{\mathbf{Z}}_t^\top \mathbf{Z}_t - \eta \mathbf{\Gamma} \mathbf{Z}_t + \eta \underline{\mathbf{\Xi}}_t \mathbf{Z}_t, \end{aligned}$$

where $\mathbf{\Gamma}$ is defined as

$$\mathbf{\Gamma} := \begin{bmatrix} \mathbf{0}_{d \times d} & \Delta \\ \Delta^\top & \mathbf{0}_{k \times k} \end{bmatrix}.$$

Then, by the triangle inequality, with probability $1 - 2C \exp(-\epsilon^2 N)$ for a universal constant $C > 0$, we have

$$\begin{aligned} \|\mathbf{Z}_{t+1}\|_{op} &\leq \|(\mathbf{I}_{2d} - \eta \mathbf{Z}_t \mathbf{Z}_t^\top) \mathbf{Z}_t\|_{op} + \eta \|\underline{\mathbf{Z}}_t \underline{\mathbf{Z}}_t^\top \mathbf{Z}_t\|_{op} + \eta \|\mathbf{\Gamma} \mathbf{Z}_t\|_{op} + \eta \|\underline{\mathbf{\Xi}}_t \mathbf{Z}_t\|_{op} \\ &\leq \left(1 - \eta \|\mathbf{Z}_t\|_{op}^2\right) \|\mathbf{Z}_t\|_{op} \quad \text{[by simultaneous diagonalization]} \\ &\quad + \eta \|\underline{\mathbf{Z}}_t \underline{\mathbf{Z}}_t^\top \mathbf{Z}_t\|_{op} + \eta \|\mathbf{\Gamma} \mathbf{Z}_t\|_{op} + \eta \|\underline{\mathbf{\Xi}}_t \mathbf{Z}_t\|_{op} \\ &\leq \left(1 - \eta \|\mathbf{Z}_t\|_{op}^2\right) \|\mathbf{Z}_t\|_{op} + \eta \|\underline{\mathbf{Z}}_t^\top \mathbf{Z}_t\|_{op} \|\mathbf{Z}_t\|_{op} + \eta \lambda_1^* \|\mathbf{Z}_t\|_{op} + \eta \epsilon \|\mathbf{R}_t\|_{op} \|\mathbf{Z}_t\|_{op}, \end{aligned}$$

where the last inequality follows from the fact that

$$\begin{aligned}\|\underline{\mathbf{Z}}_t\|_{op} &= \|\mathbf{Z}_t\|_{op}, \\ \|\mathbf{\Gamma}\|_{op} &= \lambda_1^*, \\ \|\mathbf{\Xi}_t\|_{op} &= \left\| \left(\tilde{\mathbf{\Sigma}} - \mathbf{I}_d \right) \mathbf{R}_t \right\|_{op} \leq \epsilon \|\mathbf{R}_t\|_{op}, \quad \text{w.h.p. } 1 - 2C \exp(-\epsilon^2 N).\end{aligned}$$

Using the assumption

$$\left\| \mathbf{Z}_t^\top \mathbf{Z}_t \right\|_{op} + \epsilon \|\mathbf{R}_t\|_{op} = \left\| \mathbf{A}_t^\top \mathbf{A}_t - \mathbf{B}_t^\top \mathbf{B}_t \right\|_{op} + \epsilon \|\mathbf{R}_t\|_{op} \leq \lambda_1^*,$$

then $\|\mathbf{Z}_{t+1}\|_{op}$ can be further bounded by

$$\|\mathbf{Z}_{t+1}\|_{op} \leq \left(1 - \eta \|\mathbf{Z}_t\|_{op}^2 + 2\eta\lambda_1^* \right) \|\mathbf{Z}_t\|_{op}. \quad (37)$$

Denote $x = \|\mathbf{Z}_t\|_{op}$ and $f(x) = (1 - \eta x^2 + 2\eta\lambda_1^*)x$, we have $f'(x) = 1 - 2\eta\lambda_1^* - 3\eta x^2$ and $f''(x) = -6\eta x$. Then, we know $f'(x^*) = 0$ for $x > 0$ attained at $x^* = \sqrt{\frac{1+2\eta\lambda_1^*}{3\eta}} = \sqrt{\frac{1}{3\eta} + \frac{2}{3}\lambda_1^*}$. As we pick $\eta \leq \frac{1}{10\lambda_1^*}$, then $x^* \geq 2\sqrt{\lambda_1^*}$, which implies the maximum of $f(x)$ attained at $x^* = 2\sqrt{\lambda_1^*}$ over $x \in [0, 2\lambda_1^*]$ since $\|\mathbf{Z}_t\|_{op} \leq 2\sqrt{\lambda_1^*}$ and

$$f(2\sqrt{\lambda_1^*}) = 2(1 - 4\eta\lambda_1^* + 2\eta\lambda_1^*)\sqrt{\lambda_1^*} = 2\sqrt{\lambda_1^*} - 4\eta\lambda_1^* \leq 2\sqrt{\lambda_1^*},$$

which directly implies $\|\mathbf{Z}_{t+1}\|_{op} \leq 2\sqrt{\lambda_1^*}$. By consequence, $\|\mathbf{A}_{t+1}\|_{op}, \|\mathbf{B}_{t+1}\|_{op} \leq 2\sqrt{\lambda_1^*}$ if $\|\mathbf{A}_t\|_{op}, \|\mathbf{B}_t\|_{op} \leq 2\sqrt{\lambda_1^*}$, since \mathbf{A}_{t+1} and \mathbf{B}_{t+1} are sub-matrices of \mathbf{Z}_{t+1} . \square

Based on the above results, we are ready to present the intermediate results on \mathcal{M}_t , \mathbf{A}_t , \mathbf{B}_t , and $\left\| \mathbf{A}_t^{U\perp} \right\|_{op}$.

Lemma C.16. *Under assumptions in Section 2.1 for the linear setting, with spectral initialization (Spectral-init), we take ϵ in data concentration as*

$$\epsilon \leq \min \left\{ \frac{1}{2\kappa}, \frac{\lambda_{r^*}^*}{32\kappa(32\lambda_1^* + 128\kappa^2)} \right\},$$

and set the step-size as

$$\eta \leq \min \left\{ \frac{1}{128\kappa\lambda_1^*}, \frac{(1 - \epsilon/\kappa)}{1152\lambda_1^*} \right\},$$

then with probability at least with probability $1 - 2C \exp(-\epsilon^2 N)$ for a universal constant $C > 0$, we have that $\forall t \geq 0$

$$\mathcal{M}_t \leq \frac{\lambda_{r^*}^*}{2} \quad (38)$$

$$\max \left\{ \|\mathbf{A}_t\|_{op}, \|\mathbf{B}_t\|_{op} \right\} \leq 2\sqrt{\lambda_1^*}, \quad (39)$$

$$\lambda_{r^*}^*(\mathbf{A}_t), \lambda_{r^*}^*(\mathbf{B}_t) \geq \frac{\sqrt{\lambda_{r^*}^*}}{4\sqrt{\kappa}}, \quad (40)$$

$$\left\| \mathbf{A}_t^{U\perp} \right\|_{op} \leq \frac{32\kappa\epsilon\sqrt{\lambda_1^*}}{\lambda_{r^*}^*}. \quad (41)$$

Also, we can obtain

$$\mathcal{M}_{t+1} \leq \left(1 - \eta \frac{\lambda_{r^*}^*}{64\kappa} \right) \mathcal{M}_t. \quad (42)$$

Proof. Inspired by the matrix sensing technique from Xiong et al. (2024), we develop an inductive approach to prove the claims on our settings. First, at $t = 0$, Eq. (38)-Eq. (41) can be adopted from Lemma C.11. Next, we assume Eq. (38)-Eq. (41) hold at $t \geq 1$, recall Eq. (34), by the triangle inequality, we have

$$\begin{aligned}
 \left\| \mathbf{A}_{t+1}^{U\perp} \right\|_{op} &\leq (1 - \eta \lambda_{r^*}^2 (\mathbf{B}_t^V)) \left\| \mathbf{A}_t^{U\perp} \right\|_{op} + \eta \epsilon \left\| \mathbf{R}_t \right\|_{op} \left\| \mathbf{B}_t^V \right\|_{op} \\
 &\leq (1 - \eta \lambda_{r^*}^2 (\mathbf{B}_t^V)) \left\| \mathbf{A}_t^{U\perp} \right\|_{op} + 4\eta \epsilon \mathcal{M}_t \sqrt{\lambda_1^*} \\
 &\leq \left(1 - \eta \frac{(\lambda_{r^*}^*)^2}{16\kappa} \right) \left\| \mathbf{A}_t^{U\perp} \right\|_{op} + 2\eta \epsilon \lambda_{r^*}^* \sqrt{\lambda_1^*} \\
 &\leq \left(1 - \eta \frac{(\lambda_{r^*}^*)^2}{16\kappa} \right) \frac{32\kappa \epsilon \sqrt{\lambda_1^*}}{\lambda_{r^*}^*} + 2\eta \epsilon \lambda_{r^*}^* \sqrt{\lambda_1^*} \\
 &\leq \frac{32\kappa \epsilon \sqrt{\lambda_1^*}}{\lambda_{r^*}^*},
 \end{aligned}$$

which proves the Eq. (41) at $t + 1$. Next, by Lemma C.14, we have

$$\begin{aligned}
 \left\| \mathbf{R}_{t+1}^* \right\|_{op} &\leq \left(1 - \eta \frac{\lambda_{r^*}^*}{8\kappa} \right) \mathcal{M}_t \\
 &\quad + 8\eta \epsilon \lambda_1^* \mathcal{M}_t + 2\eta^2 \lambda_1^* \lambda_{r^*}^* \mathcal{M}_t + 4\eta^2 \epsilon \lambda_1^* \lambda_{r^*}^* \mathcal{M}_t + 64\eta \epsilon \kappa^2 \mathcal{M}_t + 32\eta^2 \kappa^2 \epsilon \lambda_{r^*}^* \mathcal{M}_t + 128\eta^2 \epsilon^3 \kappa^2 \lambda_{r^*}^* \mathcal{M}_t \\
 &\quad + 64\eta^2 \epsilon^2 \kappa^2 \lambda_{r^*}^* \mathcal{M}_t + 8\eta \epsilon \lambda_1^* \mathcal{M}_t + 8\eta^2 \epsilon \lambda_1^* \mathcal{M}_t + 8\eta^2 \epsilon^2 \lambda_1^* \lambda_{r^*}^* \mathcal{M}_t + 128\eta \epsilon^2 \kappa^2 \mathcal{M}_t + 64\eta^2 \epsilon^2 \kappa^2 \lambda_{r^*}^* \mathcal{M}_t \\
 &= \left(1 - \eta \frac{\lambda_{r^*}^*}{8\kappa} \right) \mathcal{M}_t \\
 &\quad + \eta \left\{ 16\epsilon \lambda_1^* + 64\epsilon \kappa^2 + 2\eta \lambda_1^* \lambda_{r^*}^* + \eta \epsilon (4\lambda_1^* \lambda_{r^*}^* + 32\kappa^2 \lambda_{r^*}^* + 8\lambda_1^*) + 128\epsilon^2 \kappa^2 \right. \\
 &\quad \left. + \eta (128\eta \epsilon^2 \kappa^2 \lambda_{r^*}^* + 8\eta \epsilon^2 \lambda_1^* \lambda_{r^*}^*) + 128\eta \epsilon^3 \kappa^2 \lambda_{r^*}^* \right\} \mathcal{M}_t \\
 &\leq \left(1 - \eta \frac{\lambda_{r^*}^*}{8\kappa} \right) \mathcal{M}_t + 2\eta \left(16\epsilon \lambda_1^* + 64\epsilon \kappa^2 + 2\eta \lambda_1^* \lambda_{r^*}^* \right) \mathcal{M}_t && \text{[due to the order dominance]} \\
 &\leq \left(1 - \eta \frac{\lambda_{r^*}^*}{16\kappa} \right) \mathcal{M}_t + 2\eta \left(16\epsilon \lambda_1^* + 64\epsilon \kappa^2 \right) \mathcal{M}_t && \left[\text{by } \eta \leq \frac{1}{64\kappa \lambda_1^*} \right] \\
 &\leq \left(1 - \eta \frac{\lambda_{r^*}^*}{32\kappa} \right) \mathcal{M}_t, && \left[\text{by } \epsilon \leq \frac{\lambda_{r^*}^*}{16\kappa(32\lambda_1^* + 128\kappa^2)} \right]
 \end{aligned}$$

where the order dominance from the second inequality follows from the fact that η and ϵ are sufficiently small constant such that the terms in $\mathcal{O}(\eta \epsilon)$, $\mathcal{O}(\epsilon^2)$, $\mathcal{O}(\eta^2 \epsilon^2)$, $\mathcal{O}(\eta \epsilon^3)$ are significantly smaller the terms in $\mathcal{O}(\eta)$ and $\mathcal{O}(\epsilon)$.

Similarly, we can obtain

$$\begin{aligned}
 \|\mathbf{R}_{t+1}^\perp\|_{op} &\leq \left(1 - \eta \frac{\lambda_{r^*}^*}{16\kappa}\right) \mathcal{M}_t && \left[\text{since } \lambda_{\min}(\mathbf{A}_t^{U^\perp}) \geq 0\right] \\
 &+ \eta \left\{ 8\epsilon\lambda_1^* + 2\eta\lambda_1^*\lambda_{r^*}^* + 4\eta\epsilon\lambda_1^*\lambda_{r^*}^* + 64\epsilon\kappa^2 + 32\eta\epsilon\kappa^2\lambda_{r^*}^* + 64\eta\epsilon\kappa^2\lambda_{r^*}^* + 128\epsilon^2\kappa^2 \right. \\
 &\quad \left. + 8\eta\epsilon\lambda_1^* + 8\eta\epsilon^2\lambda_1^*\lambda_{r^*}^* + 2048\epsilon^3\frac{\kappa^3}{\lambda_{r^*}^*} + 64\eta\epsilon^2\kappa^2 + 128\eta\epsilon^3\kappa^2 \right\} \mathcal{M}_t \\
 &\leq \left(1 - \eta \frac{\lambda_{r^*}^*}{16\kappa}\right) \mathcal{M}_t + 2\eta \left\{ 8\epsilon\lambda_1^* + 2\eta\lambda_1^*\lambda_{r^*}^* + 64\epsilon\kappa^2 \right\} \mathcal{M}_t && [\text{due to the order dominance}] \\
 &\leq \left(1 - \eta \frac{\lambda_{r^*}^*}{32\kappa}\right) \mathcal{M}_t + 2\eta \left\{ 8\epsilon\lambda_1^* + 64\epsilon\kappa^2 \right\} \mathcal{M}_t && \left[\text{by } \eta \leq \frac{1}{128\kappa\lambda_1^*}\right] \\
 &\leq \left(1 - \eta \frac{\lambda_{r^*}^*}{64\kappa}\right) \mathcal{M}_t, && \left[\text{by } \epsilon \leq \frac{\lambda_{r^*}^*}{32\kappa(32\lambda_1^* + 128\kappa^2)}\right]
 \end{aligned}$$

which proves the Eq. (38) at $t + 1$.

Therefore, we can conclude that

$$\mathcal{M}_{t+1} \leq \left(1 - \eta \frac{\lambda_{r^*}^*}{64\kappa}\right) \mathcal{M}_t.$$

Next, assume Eq. (38)-Eq. (41) hold at $t \geq 1$, we have

$$\begin{aligned}
 (\mathbf{A}_{t+1}^\top \mathbf{A}_{t+1} - \mathbf{B}_{t+1} \mathbf{B}_{t+1}^\top) - (\mathbf{A}_t^\top \mathbf{A}_t - \mathbf{B}_t \mathbf{B}_t^\top) &= \eta^2 \mathbf{B}_t (\mathbf{A}_t \mathbf{B}_t - \Delta)^\top \widehat{\Sigma} \widehat{\Sigma} (\mathbf{A}_t \mathbf{B}_t - \Delta) \mathbf{B}_t^\top \\
 &\quad + \eta^2 \mathbf{A}_t^\top \widehat{\Sigma} (\mathbf{A}_t \mathbf{B}_t - \Delta) (\mathbf{A}_t \mathbf{B}_t - \Delta)^\top \widehat{\Sigma} \mathbf{A}_t.
 \end{aligned}$$

Accordingly, we can derive

$$\begin{aligned}
 \|(\mathbf{A}_{t+1}^\top \mathbf{A}_{t+1} - \mathbf{B}_{t+1} \mathbf{B}_{t+1}^\top) - (\mathbf{A}_0^\top \mathbf{A}_0 - \mathbf{B}_0 \mathbf{B}_0^\top)\|_{op} &= \sum_{i=1}^{t+1} \|(\mathbf{A}_i^\top \mathbf{A}_i - \mathbf{B}_i \mathbf{B}_i^\top) - (\mathbf{A}_{i-1}^\top \mathbf{A}_{i-1} - \mathbf{B}_{i-1} \mathbf{B}_{i-1}^\top)\|_{op} \\
 &= \sum_{i=1}^{t+1} 2\eta^2 \|\widehat{\Sigma}\|_{op}^2 \|\mathbf{R}_{i-1}\|_{op}^2 \max\{\|\mathbf{A}_{i-1}\|_{op}^2, \|\mathbf{B}_{i-1}\|_{op}^2\} \\
 &= \sum_{i=1}^{t+1} 72\eta^2 \mathcal{M}_{i-1}^2 \lambda_1^* && [\text{by Eq. (29)}] \\
 &\leq \sum_{i=1}^{t+1} 18\eta^2 \left(1 - \eta \frac{\lambda_{r^*}^*}{64\kappa}\right)^{2(i-1)} (\lambda_{r^*}^*)^2 \lambda_1^* \\
 &\leq 18\eta^2 (\lambda_{r^*}^*)^2 \lambda_1^* \sum_{i=0}^{\infty} \left(1 - \eta \frac{\lambda_{r^*}^*}{64\kappa}\right)^{2i} \\
 &\leq 18\eta^2 (\lambda_{r^*}^*)^2 \lambda_1^* \frac{64\kappa}{\eta\lambda_{r^*}^*} \\
 &= 1152\eta\lambda_1^*\lambda_{r^*}^*\kappa \\
 &\leq (1 - \epsilon/\kappa)\lambda_1^*. && \left[\text{by } \eta \leq \frac{(1-\epsilon/\kappa)}{1152\lambda_1^*}\right]
 \end{aligned}$$

Since $\|(\mathbf{A}_0^\top \mathbf{A}_0 - \mathbf{B}_0 \mathbf{B}_0^\top)\|_{op} = 0$ due to the spectral initialization (Spectral-init), by triangle inequality, $\|(\mathbf{A}_{t+1}^\top \mathbf{A}_{t+1} - \mathbf{B}_{t+1} \mathbf{B}_{t+1}^\top)\|_{op} \leq (1 - \epsilon/\kappa)\lambda_1^*$. Next, by Lemma C.15, we can obtain

$$\|\mathbf{A}_{t+1}\|_{op} \leq 2\sqrt{\lambda_1^*}, \quad \|\mathbf{B}_{t+1}\|_{op} \leq 2\sqrt{\lambda_1^*},$$

which proves the Eq. (39) at $t + 1$. Lastly, assume Eq. (38)-Eq. (41) hold at $t \geq 1$, by Weyl's inequality, combine with $\mathcal{M}_{t+1} \leq \frac{\lambda_{r^*}^*}{2}$, we have

$$\frac{\lambda_{r^*}^*}{2} \geq \|\mathbf{A}_{t+1}^U \mathbf{B}_{t+1}^V - \mathbf{S}^*\|_{op} \geq \lambda_{r^*}^* - \lambda_{r^*}(\mathbf{A}_{t+1}^U \mathbf{B}_{t+1}^V) \Rightarrow \lambda_{r^*}(\mathbf{A}_{t+1}^U \mathbf{B}_{t+1}^V) \geq \frac{\lambda_{r^*}^*}{2}.$$

Again by Weyl's inequality and the Eq. (39) at time $t + 1$ we can get

$$2\sqrt{\lambda_1^*} \cdot \lambda_{r^*}(\mathbf{B}_{t+1}^V) \geq \lambda_1(\mathbf{A}_{t+1}^U) \lambda_{r^*}(\mathbf{B}_{t+1}^V) \geq \lambda_{r^*}(\mathbf{A}_{t+1}^U \mathbf{B}_{t+1}^V) \geq \frac{\lambda_{r^*}^*}{2} \Rightarrow \lambda_{r^*}(\mathbf{B}_{t+1}^V) \geq \frac{\sqrt{\lambda_{r^*}^*}}{4\sqrt{\kappa}}.$$

Besides, $\lambda_{r^*}(\mathbf{A}_{t+1}^U)$ follows similar derivation. We prove all the claims. \square

Theorem C.17. *Under assumptions in Section 2.1 for the linear setting, with spectral initialization (Spectral-init), we take ϵ in data concentration as*

$$\epsilon \leq \min \left\{ \frac{1}{2\kappa}, \frac{\lambda_{r^*}^*}{32\kappa(32\lambda_1^* + 128\kappa^2)} \right\},$$

and set the step-size as

$$\eta \leq \min \left\{ \frac{1}{128\kappa\lambda_1^*}, \frac{(1 - \epsilon/\kappa)}{1152\lambda_1^*} \right\}, \quad (43)$$

then with probability at least with probability $1 - 2C \exp(-\epsilon^2 N)$ for a universal constant $C > 0$, we have that $\forall t \geq 0$

$$\|\mathbf{A}_t \mathbf{B}_t - \Delta\|_F \leq \sqrt{2r^*} \left(1 - \eta \frac{\lambda_{r^*}^*}{64\kappa} \right)^t \cdot \lambda_{r^*}^*.$$

Proof. By Lemma C.16, with probability at least with probability $1 - 2C \exp(-\epsilon^2 N)$ for a universal constant $C > 0$, we can obtain the linear convergence of generalization risk

$$\begin{aligned} \|\mathbf{A}_t \mathbf{B}_t - \Delta\|_F &\leq \sqrt{2r^*} \|\mathbf{A}_t \mathbf{B}_t - \Delta\|_{op} && [\text{Rank}(\mathbf{A}_t \mathbf{B}_t) = r^* \text{ by Lemma C.12 and Rank}(\Delta) = r^*] \\ &\leq \sqrt{2r^*} \left(1 - \eta \frac{\lambda_{r^*}^*}{64\kappa} \right)^t \cdot \lambda_{r^*}^*, \end{aligned}$$

which is independent of the choice of LoRA rank r if $r \geq r^*$. \square

Remark: The above convergence rate is independent of the choice of LoRA rank r if $r \geq r^*$. It achieves an ϵ -risk in $\mathcal{O}(\kappa^3 \ln(1/\epsilon))$ iterations. The linear convergence rate heavily depends on κ .

C.3. Preconditioned Gradient Descent under Spectral Initialization

The convergence rate in Theorem C.17 will become slow if the downstream feature shift Δ is ill-conditioned (i.e., κ is extremely large). This motivates us to add preconditioners, which is a key technique to accelerate convergence in matrix factorization/sensing (Tong et al., 2021; Zhang et al., 2021; 2023; Jia et al., 2024). The preconditioners are derived from a Riemannian metric in Mishra et al. (2012) which originally are formulated as $(\mathbf{A}\mathbf{A}^\top)^{-1}$ and $(\mathbf{B}^\top \mathbf{B})^{-1}$. Also, there is an efficient variant, i.e. $(\mathbf{A}\mathbf{A}^\top + \lambda \mathbf{I}_r)^{-1}$ and $(\mathbf{B}^\top \mathbf{B} + \lambda \mathbf{I}_r)^{-1}$, which commonly used in practice for numerical stability via the preconditioning parameter $\lambda \geq 0$.

In the over-ranked setting ($r > r^*$), $\mathbf{B}_t \mathbf{B}_t^\top$ and $\mathbf{A}_t^\top \mathbf{A}_t$ are not necessarily invertible. Hence we add the following preconditioners to vanilla GD (4)

$$\begin{aligned} \mathbf{A}_{t+1} &= \mathbf{A}_t - \frac{\eta}{N} \widetilde{\mathbf{X}}^\top \left(\widetilde{\mathbf{X}} (\mathbf{W}^\natural + \mathbf{A}_t \mathbf{B}_t) - \widetilde{\mathbf{Y}} \right) \mathbf{B}_t^\top (\mathbf{B}_t \mathbf{B}_t^\top)^\dagger, \\ \mathbf{B}_{t+1} &= \mathbf{B}_t - \frac{\eta}{N} (\mathbf{A}_t^\top \mathbf{A}_t)^\dagger \mathbf{A}_t^\top \widetilde{\mathbf{X}}^\top \left(\widetilde{\mathbf{X}} (\mathbf{W}^\natural + \mathbf{A}_t \mathbf{B}_t) - \widetilde{\mathbf{Y}} \right), \end{aligned} \quad (\text{Prec-GD})$$

where \mathbf{M}^\dagger denotes the pseudo-inverse of a matrix \mathbf{M} . Such modified preconditioners are also considered in Li et al. (2025).

In the following proofs, we will prove that the LoRA fine-tuning can achieve faster linear convergence which is independent of condition number κ under (Spectral-init) and (Prec-GD). Similar to Lemma C.12, the dynamics of \mathbf{B}_t are still limited to the r^* -dimensional singular subspace \mathbf{V} of Δ under (Spectral-init). We can verify this fact by the following lemma.

Lemma C.18. *For any natural number $t \geq 0$, under assumptions in Section 2.1 for the linear setting, with (Spectral-init) and (Prec-GD), we have*

$$\mathbf{B}_t \mathbf{V}_\perp = \mathbf{0}_{r \times (k-r^*)}.$$

Proof. For $t = 0$, recall the SVD of \mathbf{G}^\natural , i.e. $\tilde{\mathbf{U}}_{\mathbf{G}^\natural} \tilde{\mathbf{S}}_{\mathbf{G}^\natural} \tilde{\mathbf{V}}_{\mathbf{G}^\natural}^\top$ in Eq. (28), we have

$$\mathbf{B}_0 \mathbf{V}_\perp = \left[\tilde{\mathbf{S}}_{\mathbf{G}^\natural}^{-1/2} \right]_{[1:r]} \left[\tilde{\mathbf{U}}_{\mathbf{G}^\natural}^\top \right]_{[:,1:r]} \mathbf{G}^\natural \mathbf{V}_\perp = \left[\tilde{\mathbf{S}}_{\mathbf{G}^\natural}^{-1/2} \right]_{[1:r]} \left[\tilde{\mathbf{U}}_{\mathbf{G}^\natural}^\top \right]_{[:,1:r]} \hat{\Sigma} \Delta \mathbf{V}_\perp = \mathbf{0}_{r \times (k-r^*)}.$$

Assume $\mathbf{B}_t \mathbf{V}_\perp = \mathbf{0}_{d \times (d-r^*)}$ holds for any natural number $t \geq 1$, then

$$\begin{aligned} \mathbf{B}_{t+1} \mathbf{V}_\perp &= \mathbf{B}_t \mathbf{V}_\perp - \frac{\eta}{N} (\mathbf{A}_t^\top \mathbf{A}_t)^\dagger \mathbf{A}_t^\top \tilde{\mathbf{X}}^\top \left(\tilde{\mathbf{X}} (\mathbf{W}^\natural + \mathbf{A}_t \mathbf{B}_t) - \tilde{\mathbf{Y}} \right) \mathbf{V}_\perp \\ &= \mathbf{B}_t \mathbf{V}_\perp - \eta (\mathbf{A}_t^\top \mathbf{A}_t)^\dagger \mathbf{A}_t^\top \hat{\Sigma} (\mathbf{A}_t \mathbf{B}_t - \Delta) \mathbf{V}_\perp \\ &= \mathbf{0}_{r \times (k-r^*)}, \end{aligned} \quad \text{[by our inductive hypothesis]}$$

which proves the claim. \square

We can re-formulate (Prec-GD) to be

$$\mathbf{A}_{t+1} = \mathbf{A}_t - \eta \hat{\Sigma} (\mathbf{A}_t \mathbf{B}_t - \Delta) (\mathbf{B}_t)^\top (\mathbf{B}_t \mathbf{B}_t^\top)^\dagger, \quad (44)$$

$$\mathbf{B}_{t+1} = \mathbf{B}_t - \eta (\mathbf{A}_t^\top \mathbf{A}_t)^\dagger \mathbf{A}_t^\top \hat{\Sigma} (\mathbf{A}_t \mathbf{B}_t - \Delta). \quad (45)$$

Before we start our main proofs, we first define the following notations

- SVD of product matrix $\mathbf{A}_t \mathbf{B}_t := \mathcal{U}_t \mathcal{S}_t \mathcal{V}_t^\top$, where $\mathcal{U}_t \in \mathbb{R}^{d \times r^*}$, $\mathcal{S}_t \in \mathbb{R}^{r^* \times r^*}$, and $\mathcal{V}_t \in \mathbb{R}^{k \times r^*}$. Notice that here we employ rank- r^* SVD of $\mathbf{A}_t \mathbf{B}_t$ since $\text{Rank}(\mathbf{A}_t \mathbf{B}_t) \leq r^*$ due to Lemma C.18 and $\lambda_{r^*}(\mathbf{A}_t \mathbf{B}_t) > 0$ strictly which we will obtain from Theorem C.21.
- The left compact singular matrix of \mathbf{A}_t as $\mathbf{U}_{\mathbf{A}_t} \in \mathbb{R}^{d \times r}$.
- The right compact singular matrix of \mathbf{B}_t as $\mathbf{V}_{\mathbf{B}_t} \in \mathbb{R}^{k \times r^*}$. Notice that here we take the top- r^* right singular subspace of \mathbf{B}_t due to Lemma C.18.

By the pseudo inverse theorem and Jia et al. (2024, Lemma 14), we can obtain

$$\mathbf{A}_t (\mathbf{A}_t^\top \mathbf{A}_t)^\dagger \mathbf{A}_t^\top = \mathbf{U}_{\mathbf{A}_t} \mathbf{U}_{\mathbf{A}_t}^\top, \quad (46)$$

$$(\mathbf{B}_t)^\top (\mathbf{B}_t \mathbf{B}_t^\top)^\dagger \mathbf{B}_t = \mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^\top. \quad (47)$$

$$(\mathbf{B}_t)^\top (\mathbf{B}_t \mathbf{B}_t^\top)^\dagger (\mathbf{A}_t^\top \mathbf{A}_t)^\dagger \mathbf{A}_t^\top = \mathcal{V}_t \mathcal{S}_t^{-1} \mathcal{U}_t^\top. \quad (48)$$

Lemma C.19. *Denote $\mathbf{R}_t := \mathbf{A}_t \mathbf{B}_t - \Delta$, $\Xi := \hat{\Sigma} - \mathbf{I}_d$, under assumptions in Section 2.1 for the linear setting, with (Prec-GD), then we have*

$$\mathbf{R}_{t+1} = \mathbf{R}_t - \eta \mathbf{U}_{\mathbf{A}_t} \mathbf{U}_{\mathbf{A}_t}^\top \mathbf{R}_t - \eta \mathbf{R}_t \mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^\top - \eta \mathbf{U}_{\mathbf{A}_t} \mathbf{U}_{\mathbf{A}_t}^\top \Xi \mathbf{R}_t - \eta \Xi \mathbf{R}_t \mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^\top + \eta^2 \hat{\Sigma} \mathbf{R}_t \mathcal{V}_t \mathcal{S}_t^{-1} \mathcal{U}_t^\top \hat{\Sigma} \mathbf{R}_t.$$

Proof. With Eq. (44) and Eq. (45), we can construct

$$\begin{aligned}
 \mathbf{R}_{t+1} &= \mathbf{A}_{t+1} \mathbf{B}_{t+1} - \Delta \\
 &= \mathbf{A}_t \mathbf{B}_t - \Delta \\
 &\quad - \eta \mathbf{A}_t (\mathbf{A}_t^\top \mathbf{A}_t)^\dagger \mathbf{A}_t^\top \widehat{\Sigma} (\mathbf{A}_t \mathbf{B}_t - \Delta) \\
 &\quad - \eta \widehat{\Sigma} (\mathbf{A}_t \mathbf{B}_t - \Delta) (\mathbf{B}_t)^\top (\mathbf{B}_t \mathbf{B}_t^\top)^\dagger \mathbf{B}_t \\
 &\quad + \eta^2 \widehat{\Sigma} (\mathbf{A}_t \mathbf{B}_t - \Delta) (\mathbf{B}_t)^\top (\mathbf{B}_t \mathbf{B}_t^\top)^\dagger (\mathbf{A}_t^\top \mathbf{A}_t)^\dagger \mathbf{A}_t^\top \widehat{\Sigma} (\mathbf{A}_t \mathbf{B}_t - \Delta) \\
 &= \mathbf{R}_t - \eta \mathbf{U}_{\mathbf{A}_t} \mathbf{U}_{\mathbf{A}_t}^\top \widehat{\Sigma} \mathbf{R}_t - \eta \widehat{\Sigma} \mathbf{R}_t \mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^\top \quad [\text{by Eq. (46) and Eq. (47)}] \\
 &\quad + \eta^2 \widehat{\Sigma} \mathbf{R}_t \mathcal{V}_t \mathcal{S}_t^{-1} \mathcal{U}_t^\top \widehat{\Sigma} \mathbf{R}_t \quad [\text{by Eq. (48)}] \\
 &= \mathbf{R}_t - \eta \mathbf{U}_{\mathbf{A}_t} \mathbf{U}_{\mathbf{A}_t}^\top \mathbf{R}_t - \eta \mathbf{R}_t \mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^\top - \eta \mathbf{U}_{\mathbf{A}_t} \mathbf{U}_{\mathbf{A}_t}^\top \Xi \mathbf{R}_t - \eta \Xi \mathbf{R}_t \mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^\top + \eta^2 \widehat{\Sigma} \mathbf{R}_t \mathcal{V}_t \mathcal{S}_t^{-1} \mathcal{U}_t^\top \widehat{\Sigma} \mathbf{R}_t,
 \end{aligned}$$

which proves the claim. \square

In the next, we aim to estimate the signal part $\mathbf{R}_t - \eta \mathbf{U}_{\mathbf{A}_t} \mathbf{U}_{\mathbf{A}_t}^\top \mathbf{R}_t - \eta \mathbf{R}_t \mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^\top$.

Lemma C.20. Recall $\mathbf{R}_t := \mathbf{A}_t \mathbf{B}_t - \Delta$, under assumptions in Section 2.1 for the linear setting, with (Prec-GD), then

$$\|\mathbf{R}_t - \eta \mathbf{U}_{\mathbf{A}_t} \mathbf{U}_{\mathbf{A}_t}^\top \mathbf{R}_t - \eta \mathbf{R}_t \mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^\top\|_F \leq (1 - \eta) \|\mathbf{R}_t\|_F.$$

Proof.

$$\begin{aligned}
 &\|\mathbf{R}_t - \eta \mathbf{U}_{\mathbf{A}_t} \mathbf{U}_{\mathbf{A}_t}^\top \mathbf{R}_t - \eta \mathbf{R}_t \mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^\top\|_F \\
 &= \|\mathbf{R}_t (\mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^\top + \mathbf{I}_k - \mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^\top) - \eta \mathbf{U}_{\mathbf{A}_t} \mathbf{U}_{\mathbf{A}_t}^\top \mathbf{R}_t (\mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^\top + \mathbf{I}_k - \mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^\top) - \eta \mathbf{R}_t \mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^\top\|_F \\
 &= \|\mathbf{R}_t \mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^\top - \eta \mathbf{U}_{\mathbf{A}_t} \mathbf{U}_{\mathbf{A}_t}^\top \mathbf{R}_t \mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^\top - \eta \mathbf{R}_t \mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^\top\|_F \quad [\text{since } \mathbf{R}_t (\mathbf{I}_k - \mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^\top) = \mathbf{0} \text{ by Lemma C.18}] \\
 &= \|(\mathbf{I}_d - \eta (\mathbf{I}_d + \mathbf{U}_{\mathbf{A}_t} \mathbf{U}_{\mathbf{A}_t}^\top)) \mathbf{R}_t \mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^\top\|_F \\
 &= \|\mathbf{I}_d - \eta (\mathbf{I}_d + \mathbf{U}_{\mathbf{A}_t} \mathbf{U}_{\mathbf{A}_t}^\top)\|_{op} \|\mathbf{R}_t \mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^\top\|_F \\
 &\leq (1 - \eta) \|\mathbf{R}_t\|_F, \quad \left[\|\mathbf{I}_d - \eta (\mathbf{I}_d + \mathbf{U}_{\mathbf{A}_t} \mathbf{U}_{\mathbf{A}_t}^\top)\|_{op} \leq 1 - \eta, \text{ since } \mathbf{U}_{\mathbf{A}_t} \mathbf{U}_{\mathbf{A}_t}^\top \text{ is a rank-}r \text{ projection matrix} \right]
 \end{aligned}$$

which concludes the proof. \square

Finally, we have the following linear convergence under (Prec-GD) and (Spectral-init).

Theorem C.21. Under assumptions in Section 2.1 for the linear setting, with (Spectral-init) and (Prec-GD), we choose

$$\epsilon \leq \min \left\{ \frac{1}{2\sqrt{r^* \kappa}}, \frac{1}{4} \right\}$$

and set $\eta \in \left(0, \frac{0.5-2\epsilon}{(1+\epsilon)^2}\right)$, then with probability at least $1 - 2C \exp(-\epsilon^2 N)$ for a universal constant $C > 0$, we have

$$\|\mathbf{A}_t \mathbf{B}_t - \Delta\|_F \leq \frac{1}{2} \left(1 - \frac{\eta}{2}\right)^t \lambda_{r^*}^*.$$

Proof. We prove it by induction. We suppose the following two inductive hypothesis

$$\lambda_{r^*}(\mathbf{A}_t \mathbf{B}_t) \geq \frac{\lambda_{r^*}^*}{2}, \quad (49)$$

$$\|\mathbf{A}_0 \mathbf{B}_0 - \Delta\|_F \leq \frac{\lambda_{r^*}^*}{2}. \quad (50)$$

Starting from $t = 0$, under (Spectral-init), with probability at least $1 - 2C \exp(-\epsilon^2 N)$ for a universal constant $C > 0$, we have

$$\begin{aligned}
 \|\mathbf{A}_0 \mathbf{B}_0 - \Delta\|_F &= \|\mathbf{G}^\natural - \Delta\|_F \\
 &= \left\| \left(\widehat{\Sigma} - \mathbf{I}_d \right) \Delta \right\|_F && \text{[by Eq. (28)]} \\
 &\leq \epsilon \|\Delta\|_F \\
 &\leq \epsilon \sqrt{r^*} \|\Delta\|_{op} && \text{[since Rank}(\Delta) = r^* \text{]} \\
 &\leq \frac{\lambda_{r^*}^*}{2}. && \text{[since } \epsilon \leq 1/2\sqrt{r^* \kappa} \text{]}
 \end{aligned}$$

Then, by Weyl's inequality, we have

$$\lambda_{r^*}(\Delta) - \lambda_{r^*}(\mathbf{A}_0 \mathbf{B}_0) \leq \|\mathbf{A}_0 \mathbf{B}_0 - \Delta\|_{op} \leq \|\mathbf{A}_0 \mathbf{B}_0 - \Delta\|_F,$$

which implies

$$\lambda_{r^*}(\mathbf{A}_0 \mathbf{B}_0) \geq \frac{\lambda_{r^*}^*}{2}. \quad (51)$$

Therefore, we verify Eq. (49) and Eq. (50) at $t = 0$. We assume Eq. (49) and Eq. (50) hold at $t = 2, 3, \dots$, then by Lemma C.19, with probability at least with probability $1 - 2C \exp(-\epsilon^2 N)$ for a universal constant $C > 0$, we have

$$\begin{aligned}
 \|\mathbf{R}_{t+1}\|_F &\leq \left\| \mathbf{R}_t - \eta \mathbf{U}_{\mathbf{A}_t} \mathbf{U}_{\mathbf{A}_t}^\top \mathbf{R}_t - \eta \mathbf{R}_t \mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^\top \right\|_F \\
 &\quad + \eta \left\| \mathbf{U}_{\mathbf{A}_t} \mathbf{U}_{\mathbf{A}_t}^\top \Xi \mathbf{R}_t \right\|_F + \eta \left\| \Xi \mathbf{R}_t \mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^\top \right\|_F + \eta^2 \left\| \widehat{\Sigma} \mathbf{R}_t \mathcal{V}_t \mathcal{S}_t^{-1} \mathcal{U}_t^\top \widehat{\Sigma} \mathbf{R}_t \right\|_F \\
 &\leq (1 - \eta) \|\mathbf{R}_t\|_F && \text{[by Lemma C.20]} \\
 &\quad + \eta \epsilon \left\| \mathbf{U}_{\mathbf{A}_t} \mathbf{U}_{\mathbf{A}_t}^\top \mathbf{R}_t \right\|_F + \eta \epsilon \left\| \mathbf{R}_t \mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^\top \right\|_F + \eta^2 (1 + \epsilon)^2 \frac{\|\mathbf{R}_t\|_F^2}{\lambda_{r^*}(\mathbf{A}_t \mathbf{B}_t)} && \text{[by } \|\Xi\|_{op} \leq \epsilon \text{]} \\
 &\leq (1 - \eta) \|\mathbf{R}_t\|_F \\
 &\quad + \eta \epsilon \|\mathbf{R}_t\|_F + \eta \epsilon \|\mathbf{R}_t\|_F + \eta^2 (1 + \epsilon)^2 \|\mathbf{R}_t\|_F && \text{[since Eq. (49) and Eq. (50) hold at } t \text{]} \\
 &= (1 - (1 - 2\epsilon)\eta + \eta^2 (1 + \epsilon)^2) \|\mathbf{R}_t\|_F \\
 &\leq \left(1 - \frac{\eta}{2}\right) \|\mathbf{R}_t\|_F. && \text{[taking } \eta \leq \frac{0.5 - 2\epsilon}{(1 + \epsilon)^2} \text{]}
 \end{aligned}$$

This implies Eq. (50) at time $t + 1$. By consequence, we can obtain Eq. (49) at time $t + 1$ again by Weyl's inequality. \square

Remark: The convergence rate is independent of the condition number of κ . The choice of stepsize η is upper bounded by $\frac{0.5 - 2\epsilon}{(1 + \epsilon)^2} \in (0, 0.5)$, which is a decreasing function of ϵ . Therefore, if the condition number κ is very large and thus ϵ is chosen as sufficiently small, then η can reach 0.5 and we still have a fast convergence rate independent of κ . This is particularly useful in practical fine-tuning tasks, where the adapted matrix can be highly ill-conditioned when its rank increases. We can empirically observe the ill-conditioned issues in real-world benchmarks, as shown in Appendix G.5 for more discussions.

D. Proofs for Nonlinear Model

We deliver the proofs for nonlinear models in Section 4 here. The problem setting and results for $\|\mathbf{A}_0 \mathbf{B}_0 - \Delta\|_F$ are presented in Appendix D.1. In Appendix D.2, we present the proofs of Theorem 4.2 as well as proofs for smoothed GD.

D.1. Problem Settings and Spectral Initialization

Recall the pre-training model from Assumption 2.1

$$f_{\text{pre}}(\mathbf{x}) = \sigma(\mathbf{x}^\top \mathbf{W}^\natural)^\top \in \mathbb{R}^k, \quad \mathbf{W}^\natural \in \mathbb{R}^{d \times k}, \quad (52)$$

and the downstream teacher weights from Assumption 2.2

$$\widetilde{\mathbf{W}}^\natural = \mathbf{W}^\natural + \Delta \in \mathbb{R}^{d \times k}, \quad \text{with } \widetilde{\mathbf{W}}^\natural := [\widetilde{\mathbf{w}}_1^\natural, \widetilde{\mathbf{w}}_2^\natural, \dots, \widetilde{\mathbf{w}}_k^\natural].$$

The empirical loss of LoRA fine-tuning is defined as

$$\widetilde{L}(\mathbf{A}_t, \mathbf{B}_t) = \frac{1}{2N} \left\| \sigma(\widetilde{\mathbf{X}}(\mathbf{W}^\natural + \mathbf{A}_t \mathbf{B}_t)) - \sigma(\widetilde{\mathbf{X}} \widetilde{\mathbf{W}}^\natural) \right\|_{\text{F}}^2.$$

Next, we can derive the empirical gradients for \mathbf{A}_t and \mathbf{B}_t respectively.

$$\begin{aligned} \nabla_{\mathbf{A}} \widetilde{L}(\mathbf{A}_t, \mathbf{B}_t) &= \frac{1}{N} \widetilde{\mathbf{X}}^\top \left[\sigma(\widetilde{\mathbf{X}}(\mathbf{W}^\natural + \mathbf{A}_t \mathbf{B}_t)) - \sigma(\widetilde{\mathbf{X}} \widetilde{\mathbf{W}}^\natural) \right] \odot \sigma'(\widetilde{\mathbf{X}}(\mathbf{W}^\natural + \mathbf{A}_t \mathbf{B}_t)) \mathbf{B}_t^\top \\ &:= \frac{1}{N} \widetilde{\mathbf{X}}^\top \left[\sigma(\widetilde{\mathbf{X}}(\mathbf{W}_t)) - \sigma(\widetilde{\mathbf{X}} \widetilde{\mathbf{W}}^\natural) \right] \odot \sigma'(\widetilde{\mathbf{X}} \mathbf{W}_t) \mathbf{B}_t^\top \quad [\text{denote } \mathbf{W}_t := \mathbf{W}^\natural + \mathbf{A}_t \mathbf{B}_t] \\ &= - \left[\underbrace{\frac{1}{N} \widetilde{\mathbf{X}}^\top \sigma(\widetilde{\mathbf{X}} \widetilde{\mathbf{W}}^\natural) \odot \sigma'(\widetilde{\mathbf{X}} \mathbf{W}_t)}_{:= \Gamma_{1,t}} - \underbrace{\frac{1}{N} \widetilde{\mathbf{X}}^\top \sigma(\widetilde{\mathbf{X}} \mathbf{W}_t) \odot \sigma'(\widetilde{\mathbf{X}} \mathbf{W}_t)}_{:= \Gamma_{2,t}} \right] \mathbf{B}_t^\top \\ &:= -\mathbf{J}_{\mathbf{W}_t} \mathbf{B}_t^\top \quad [\text{denote } \mathbf{J}_{\mathbf{W}_t} := \Gamma_{1,t} - \Gamma_{2,t}] \end{aligned} \quad (53)$$

where the matrix operator $\mathbf{J}_{\mathbf{W}} : \mathbb{R}^{d \times k} \rightarrow \mathbb{R}^{d \times k}$ is formally defined as (by denoting $\mathbf{W}_t := \mathbf{W}^\natural + \mathbf{A}_t \mathbf{B}_t$)

$$\mathbf{J}_{\mathbf{W}} : \mathbf{W} \rightarrow \frac{1}{N} \widetilde{\mathbf{X}}^\top \left[\sigma(\widetilde{\mathbf{X}} \widetilde{\mathbf{W}}^\natural) - \sigma(\widetilde{\mathbf{X}}(\mathbf{W})) \right] \odot \sigma'(\widetilde{\mathbf{X}} \mathbf{W}). \quad (54)$$

Similarly, we can compute

$$\begin{aligned} \nabla_{\mathbf{B}} \widetilde{L}(\mathbf{A}_t, \mathbf{B}_t) &= \frac{1}{N} \mathbf{A}_t^\top \widetilde{\mathbf{X}}^\top \left[\sigma(\widetilde{\mathbf{X}}(\mathbf{W}_t)) - \sigma(\widetilde{\mathbf{X}} \widetilde{\mathbf{W}}^\natural) \right] \odot \sigma'(\widetilde{\mathbf{X}} \mathbf{W}_t) \\ &= -\mathbf{A}_t^\top \mathbf{J}_{\mathbf{W}_t}. \end{aligned}$$

For full fine-tuning, we consider the following empirical loss function over $\mathbf{K} \in \mathbb{R}^{d \times k}$

$$L(\mathbf{K}) = \frac{1}{2N} \left\| \sigma(\widetilde{\mathbf{X}} \mathbf{K}) - \sigma(\widetilde{\mathbf{X}} \widetilde{\mathbf{W}}^\natural) \right\|_{\text{F}}^2.$$

The gradient w.r.t. \mathbf{K} is

$$\nabla L(\mathbf{K}) = \frac{1}{N} \widetilde{\mathbf{X}}^\top \left[\sigma(\widetilde{\mathbf{X}} \mathbf{K}) - \sigma(\widetilde{\mathbf{X}} \widetilde{\mathbf{W}}^\natural) \right] \odot \sigma'(\widetilde{\mathbf{X}} \mathbf{K})$$

Next, we can define the one-step negative gradient of full fine-tuning in the nonlinear case as

$$\begin{aligned} \mathbf{G}^\natural &:= -\nabla L(\mathbf{W}^\natural) \\ &= \frac{1}{N} \widetilde{\mathbf{X}}^\top \left[\sigma(\widetilde{\mathbf{X}} \widetilde{\mathbf{W}}^\natural) - \sigma(\widetilde{\mathbf{X}} \mathbf{W}^\natural) \right] \odot \sigma'(\widetilde{\mathbf{X}} \mathbf{W}^\natural) \\ &= \mathbf{J}_{\mathbf{W}^\natural}. \quad [\text{by definition of } \mathbf{J}_{\mathbf{W}} \text{ in Eq. (54)}] \end{aligned}$$

Additionally, we define

$$\begin{aligned} \Gamma_1^\natural &= \frac{1}{N} \widetilde{\mathbf{X}}^\top \sigma(\widetilde{\mathbf{X}} \widetilde{\mathbf{W}}^\natural) \odot \sigma'(\widetilde{\mathbf{X}} \mathbf{W}^\natural), \\ \Gamma_2^\natural &= \frac{1}{N} \widetilde{\mathbf{X}}^\top \sigma(\widetilde{\mathbf{X}} \mathbf{W}^\natural) \odot \sigma'(\widetilde{\mathbf{X}} \mathbf{W}^\natural). \end{aligned} \quad (55)$$

In this section, we aim to analyze the initial properties of low-rank adapters under (Spectral-init) in a nonlinear context. The high-level proof strategy begins with examining the spectral properties of the one-step full gradient matrix, \mathbf{G}^\natural . Unlike the linear case, the presence of nonlinearity prevents a direct analysis. To address this, we first establish the concentration of the

empirical full gradient, leveraging the fact that the empirical gradient approximates its expectation closely when the sample size is sufficiently large.

Subsequently, we utilize tools from (Brutzkus & Globerson, 2017) to derive useful properties of the expected gradients. These properties are then transferred back to the empirical gradients through concentration results. Finally, since low-rank adapters under (Spectral-init) represent the best r -rank approximation of G^\natural , we apply matrix analysis techniques to derive the desired results. Also, the concentration results in this part can serve as an important component for the later convergence analysis.

D.1.1. COMPUTATION OF FULL POPULATION GRADIENTS

First, we can simplify $\Gamma_{1,t}$ and $\Gamma_{2,t}$ which defined in Eq. (53) to be

$$\Gamma_{1,t} = \frac{1}{N} \sum_{i=1}^N \tilde{\mathbf{x}}_i \left[\sigma(\tilde{\mathbf{x}}_i^\top \tilde{\mathbf{w}}_1^\natural) \sigma'(\tilde{\mathbf{x}}_i^\top \mathbf{w}_{t,1}) \quad \dots \quad \sigma(\tilde{\mathbf{x}}_i^\top \tilde{\mathbf{w}}_k^\natural) \sigma'(\tilde{\mathbf{x}}_i^\top \mathbf{w}_{t,k}) \right],$$

and

$$\Gamma_{2,t} = \frac{1}{N} \sum_{i=1}^N \tilde{\mathbf{x}}_i \left[\sigma(\tilde{\mathbf{x}}_i^\top \mathbf{w}_{t,1}) \sigma'(\tilde{\mathbf{x}}_i^\top \mathbf{w}_{t,1}) \quad \dots \quad \sigma(\tilde{\mathbf{x}}_i^\top \mathbf{w}_{t,k}) \sigma'(\tilde{\mathbf{x}}_i^\top \mathbf{w}_{t,k}) \right],$$

where $\mathbf{w}_{t,m}$ is the m -th column of $\mathbf{W}_t := \tilde{\mathbf{W}}^\natural + \mathbf{A}_t \mathbf{B}_t$ and $\tilde{\mathbf{w}}_m^\natural$ is the m -th column of $\tilde{\mathbf{W}}^\natural$.

The following two lemmas provide the columnwise expectation of $\Gamma_{1,t}$ and $\Gamma_{2,t}$ respectively.

Lemma D.1. *Under assumptions in Section 2.1 for the nonlinear setting, for any $1 \leq m \leq k$, we have*

$$\mathbb{E}_{\tilde{\mathbf{x}}} [\tilde{\mathbf{x}} \sigma'(\tilde{\mathbf{x}}^\top \mathbf{w}_{t,m}) \sigma(\tilde{\mathbf{x}}^\top \tilde{\mathbf{w}}_m^\natural)] = \frac{1}{2\pi} \left[\frac{\|\tilde{\mathbf{w}}_m^\natural\|_2}{\|\mathbf{w}_{t,m}\|_2} \sin \theta(\mathbf{w}_{t,m}, \tilde{\mathbf{w}}_m^\natural) \mathbf{w}_{t,m} + (\pi - \theta(\mathbf{w}_{t,m}, \tilde{\mathbf{w}}_m^\natural)) \tilde{\mathbf{w}}_m^\natural \right], \quad (56)$$

and

$$\mathbb{E}_{\tilde{\mathbf{x}}} [\mathbf{x} \sigma'(\tilde{\mathbf{x}}^\top \mathbf{w}_{t,m}) \sigma(\tilde{\mathbf{x}}^\top \mathbf{w}_{t,m})] = \frac{1}{2} \mathbf{w}_{t,m}. \quad (57)$$

Proof. First, for any $1 \leq m \leq k$, we have

$$\begin{aligned} & \mathbb{E}_{\tilde{\mathbf{x}}} [\sigma'(\tilde{\mathbf{x}}^\top \mathbf{w}_{t,m}) \sigma(\tilde{\mathbf{x}}^\top \tilde{\mathbf{w}}_m^\natural) \tilde{\mathbf{x}}] \\ &= \frac{\partial}{\partial \mathbf{w}_{t,m}} \mathbb{E}_{\tilde{\mathbf{x}}} \left[\sigma(\tilde{\mathbf{x}}^\top \mathbf{w}_{t,m}) \sigma(\tilde{\mathbf{x}}^\top \tilde{\mathbf{w}}_m^\natural) \right] \\ &= \frac{1}{2\pi} \left[\frac{\|\tilde{\mathbf{w}}_m^\natural\|_2}{\|\mathbf{w}_{t,m}\|_2} \sin \theta(\mathbf{w}_{t,m}, \tilde{\mathbf{w}}_m^\natural) \mathbf{w}_{t,m} + (\pi - \theta(\mathbf{w}_{t,m}, \tilde{\mathbf{w}}_m^\natural)) \tilde{\mathbf{w}}_m^\natural \right]. \end{aligned} \quad [\text{by Lemma E.4}]$$

The proof for Eq. (57) is the same as that of Eq. (56), i.e.

$$\begin{aligned} & \mathbb{E}_{\tilde{\mathbf{x}}} [\sigma'(\tilde{\mathbf{x}}^\top \mathbf{w}_{t,m}) \sigma(\tilde{\mathbf{x}}^\top \mathbf{w}_{t,m}) \tilde{\mathbf{x}}] \\ &= \frac{1}{2\pi} \left[\frac{\|\mathbf{w}_{t,m}\|_2}{\|\mathbf{w}_{t,m}\|_2} \sin \theta(\mathbf{w}_{t,m}, \mathbf{w}_{t,m}) \mathbf{w}_{t,m} + (\pi - \theta(\mathbf{w}_{t,m}, \mathbf{w}_{t,m})) \mathbf{w}_{t,m} \right] \\ &= \frac{1}{2} \mathbf{w}_{t,m}. \end{aligned} \quad [\text{by Lemma E.4}]$$

□

Next, we can obtain the expected full gradients via the following lemma.

Lemma D.2. *Recall $\mathbf{W}_t := \mathbf{W}^\natural + \mathbf{A}_t \mathbf{B}_t$, suppose $\|\mathbf{A}_t \mathbf{B}_t - \Delta\|_F \leq \rho \lambda_{r^*}^*$ for some constant $\rho \in (0, 1)$, under assumptions in Section 2.1 for the nonlinear setting and Assumption 4.1, then it holds that*

$$-\mathbb{E}_{\tilde{\mathbf{x}}} [\mathbf{J}_{\mathbf{W}_t}] = \frac{1}{2} (\mathbf{A}_t \mathbf{B}_t - \Delta) + \Psi(t),$$

where $\Psi(t)$ is defined as

$$\Psi(t) := (\mathbf{A}_t \mathbf{B}_t - \Delta) \mathbf{D}_1(t) \mathbf{D}_3(t) + \widetilde{\mathbf{W}}^\natural \left[\mathbf{D}_1(t) - \mathbf{D}_2(t) - \mathbf{D}_1(t)(\mathbf{I}_k - \mathbf{D}_3(t)) \right],$$

with

$$\begin{aligned} \mathbf{D}_1(t) &:= \text{Diag} \left\{ \sin \left[\theta(\widetilde{\mathbf{w}}_1^\natural + \mathbf{r}_{t,1}, \widetilde{\mathbf{w}}_1^\natural) \right], \dots, \sin \left[\theta(\widetilde{\mathbf{w}}_k^\natural + \mathbf{r}_{t,k}, \widetilde{\mathbf{w}}_k^\natural) \right] \right\}, \\ \mathbf{D}_2(t) &:= \text{Diag} \left\{ \theta(\widetilde{\mathbf{w}}_1^\natural + \mathbf{r}_{t,1}, \widetilde{\mathbf{w}}_1^\natural), \dots, \theta(\widetilde{\mathbf{w}}_k^\natural + \mathbf{r}_{t,k}, \widetilde{\mathbf{w}}_k^\natural) \right\}, \\ \mathbf{D}_3(t) &:= \text{Diag} \left\{ \frac{\|\widetilde{\mathbf{w}}_1^\natural\|_2}{\|\widetilde{\mathbf{w}}_1^\natural + \mathbf{r}_{t,1}\|_2}, \dots, \frac{\|\widetilde{\mathbf{w}}_k^\natural\|_2}{\|\widetilde{\mathbf{w}}_k^\natural + \mathbf{r}_{t,k}\|_2} \right\}. \end{aligned} \quad (58)$$

Then we have the following upper bound

$$\frac{\|\Psi(t)\|_F}{\|\mathbf{A}_t \mathbf{B}_t - \Delta\|_F} \leq \mathcal{O} \left(\frac{1}{\kappa r^*} \right).$$

Proof. We first give some notations here. Let $\mathbf{w}_{t,m}$ be the m -th column of $\mathbf{W}_t \in \mathbb{R}^{d \times k}$, \mathbf{w}_m^\natural be the m -th column of $\widetilde{\mathbf{W}}^\natural$, Δ_m as the m -th of the low-rank shift Δ , $[\mathbf{A}_t \mathbf{B}_t]_m$ as the m -th column of $\mathbf{A}_t \mathbf{B}_t$.

By Lemma D.1, we can derive m -th column of $-\mathbb{E}_{\widetilde{\mathbf{x}}} [\mathbf{J}_{\mathbf{W}_t}]$ for any $m = 1, 2, \dots, k$ as

$$\begin{aligned} & \mathbb{E}_{\widetilde{\mathbf{x}}} \left[\frac{1}{N} \sum_{i=1}^N (\sigma(\widetilde{\mathbf{x}}_i^\top \mathbf{w}_{t,m}) - \sigma(\widetilde{\mathbf{x}}_i^\top \mathbf{w}_m^\natural)) \sigma'(\widetilde{\mathbf{x}}_i^\top \mathbf{w}_{t,m}) \widetilde{\mathbf{x}}_i \right] \\ &= \mathbb{E}_{\widetilde{\mathbf{x}}} \left[(\sigma(\widetilde{\mathbf{x}}^\top \mathbf{w}_{t,m}) - \sigma(\widetilde{\mathbf{x}}^\top \mathbf{w}_m^\natural)) \sigma'(\widetilde{\mathbf{x}}^\top \mathbf{w}_{t,m}) \widetilde{\mathbf{x}} \right] \\ &= \frac{1}{2} (\mathbf{w}_{t,m} - \mathbf{w}_m^\natural) - \frac{1}{2\pi} \left[\frac{\|\widetilde{\mathbf{w}}_m^\natural\|_2}{\|\mathbf{w}_{t,m}\|_2} \sin \theta(\mathbf{w}_{t,m}, \widetilde{\mathbf{w}}_m^\natural) \mathbf{w}_{t,m} - \theta(\mathbf{w}_{t,m}, \widetilde{\mathbf{w}}_m^\natural) \widetilde{\mathbf{w}}_m^\natural \right] \\ &= \frac{1}{2} ([\mathbf{A}_t \mathbf{B}_t]_m - \Delta_m) - \underbrace{\frac{1}{2\pi} \left[\frac{\|\widetilde{\mathbf{w}}_m^\natural\|_2}{\|\mathbf{w}_{t,m}\|_2} \sin \theta(\mathbf{w}_{t,m}, \widetilde{\mathbf{w}}_m^\natural) \mathbf{w}_{t,m} - \theta(\mathbf{w}_{t,m}, \widetilde{\mathbf{w}}_m^\natural) \widetilde{\mathbf{w}}_m^\natural \right]}_{\text{residual part R}}. \end{aligned} \quad (59)$$

Denote

$$\mathbf{r}_{t,m} := [\mathbf{A}_t \mathbf{B}_t]_m - \Delta_m, \quad (60)$$

then we can write

$$\mathbf{w}_{t,m} = \widetilde{\mathbf{w}}_m^\natural + [\mathbf{A}_t \mathbf{B}_t]_m - \Delta_m = \widetilde{\mathbf{w}}_m^\natural + \mathbf{r}_{t,m},$$

as a relative perturbed time-dependent vector. Next, we take it back to the residual part in Eq. (59)

$$\begin{aligned} \mathbf{R} &= \frac{\|\widetilde{\mathbf{w}}_m^\natural\|_2}{\|\widetilde{\mathbf{w}}_m^\natural + \mathbf{r}_{t,m}\|_2} \sin [\theta(\widetilde{\mathbf{w}}_m^\natural + \mathbf{r}_{t,m}, \widetilde{\mathbf{w}}_m^\natural)] (\widetilde{\mathbf{w}}_m^\natural + \mathbf{r}_{t,m}) - \theta(\widetilde{\mathbf{w}}_m^\natural + \mathbf{r}_{t,m}, \widetilde{\mathbf{w}}_m^\natural) \widetilde{\mathbf{w}}_m^\natural \\ &= \frac{\|\widetilde{\mathbf{w}}_m^\natural\|_2}{\|\widetilde{\mathbf{w}}_m^\natural + \mathbf{r}_{t,m}\|_2} \sin [\theta(\widetilde{\mathbf{w}}_m^\natural + \mathbf{r}_{t,m}, \widetilde{\mathbf{w}}_m^\natural)] \mathbf{r}_{t,m} \\ &\quad - \frac{\|\widetilde{\mathbf{w}}_m^\natural + \mathbf{r}_{t,m}\|_2 - \|\widetilde{\mathbf{w}}_m^\natural\|_2}{\|\widetilde{\mathbf{w}}_m^\natural + \mathbf{r}_{t,m}\|_2} \sin [\theta(\widetilde{\mathbf{w}}_m^\natural + \mathbf{r}_{t,m}, \widetilde{\mathbf{w}}_m^\natural)] \widetilde{\mathbf{w}}_m^\natural \\ &\quad + \sin [\theta(\widetilde{\mathbf{w}}_m^\natural + \mathbf{r}_{t,m}, \widetilde{\mathbf{w}}_m^\natural)] \widetilde{\mathbf{w}}_m^\natural - \theta(\widetilde{\mathbf{w}}_m^\natural + \mathbf{r}_{t,m}, \widetilde{\mathbf{w}}_m^\natural) \widetilde{\mathbf{w}}_m^\natural. \end{aligned} \quad (61)$$

Combining Eq. (59) and Eq. (61), we can write $-\mathbb{E}_{\widetilde{\mathbf{x}}} [\mathbf{J}_{\mathbf{W}_t}]$ in matrix form as

$$\mathbb{E}_{\widetilde{\mathbf{x}}} [\mathbf{J}_{\mathbf{W}_t}] = \frac{1}{2} (\mathbf{A}_t \mathbf{B}_t - \Delta) + (\mathbf{A}_t \mathbf{B}_t - \Delta) \mathbf{D}_1(t) \mathbf{D}_3(t) + \widetilde{\mathbf{W}}^\natural \left[\mathbf{D}_1(t) - \mathbf{D}_2(t) - \mathbf{D}_1(t)(\mathbf{I}_k - \mathbf{D}_3(t)) \right].$$

Additionally, we define

$$\begin{aligned}\Psi(t) &:= (\mathbf{A}_t \mathbf{B}_t - \Delta) \mathbf{D}_1(t) \mathbf{D}_3(t) + \widetilde{\mathbf{W}}^\natural \left[\mathbf{D}_1(t) - \mathbf{D}_2(t) - \mathbf{D}_1(t)(\mathbf{I}_k - \mathbf{D}_3(t)) \right] \\ &= \underbrace{(\mathbf{A}_t \mathbf{B}_t - \Delta) \mathbf{D}_1(t) \mathbf{D}_3(t)}_{:= \Psi_1(t)} + \underbrace{\widetilde{\mathbf{W}}^\natural (\mathbf{D}_1(t) - \mathbf{D}_2(t))}_{:= \Psi_2(t)} - \underbrace{\widetilde{\mathbf{W}}^\natural (\mathbf{D}_1(t)(\mathbf{I}_k - \mathbf{D}_3(t)))}_{:= \Psi_3(t)}.\end{aligned}\quad (62)$$

For notational simplicity, we drop time & column index and denote $\mathbf{w} := \widetilde{\mathbf{w}}_m^\natural$ and $\mathbf{r} := \mathbf{r}_{t,m}$. By condition $\|\mathbf{A}_t \mathbf{B}_t - \Delta\|_F \leq \rho \lambda_{r^*}^*$, we have

$$\|\mathbf{r}\|_2 \leq \|\mathbf{A}_t \mathbf{B}_t - \Delta\|_F \leq \rho \lambda_{r^*}^*.$$

Next, by Assumption 4.1, we can obtain

$$\frac{\|\mathbf{r}\|_2}{\|\mathbf{w}\|_2} \leq \frac{\rho \lambda_{r^*}^*}{\|\mathbf{w}\|_2} \leq \mathcal{O}\left(\frac{1}{\kappa r^*}\right). \quad (63)$$

Part I: Control the Angle $\theta(\mathbf{w} + \mathbf{r}, \mathbf{w})$

We denote $\alpha := \frac{\langle \mathbf{r}, \mathbf{w} \rangle}{\|\mathbf{w}\|_2^2}$ and $\beta := \frac{\|\mathbf{r}\|_2^2}{\|\mathbf{w}\|_2^2}$, then we can derive

$$\cos \theta(\mathbf{w} + \mathbf{r}, \mathbf{w}) = \frac{1 + \frac{\langle \mathbf{r}, \mathbf{w} \rangle}{\|\mathbf{w}\|_2^2}}{\sqrt{1 + 2\frac{\langle \mathbf{r}, \mathbf{w} \rangle}{\|\mathbf{w}\|_2^2} + \frac{\|\mathbf{r}\|_2^2}{\|\mathbf{w}\|_2^2}}} = \frac{1 + \alpha}{\sqrt{1 + 2\alpha + \beta}},$$

which can imply

$$\sin \theta(\mathbf{w} + \mathbf{r}, \mathbf{w}) = \sqrt{1 - \cos^2 \theta(\mathbf{w} + \mathbf{r}, \mathbf{w})} = \sqrt{\frac{\beta - \alpha^2}{1 + 2\alpha + \beta}} = \Theta\left(\sqrt{\beta - \alpha^2}\right) = \mathcal{O}\left(\frac{\|\mathbf{r}\|_2}{\|\mathbf{w}\|_2}\right). \quad (64)$$

By consequence, we can obtain

$$\theta(\mathbf{w} + \mathbf{r}, \mathbf{w}) = \mathcal{O}\left(\frac{\|\mathbf{r}\|_2}{\|\mathbf{w}\|_2}\right). \quad (65)$$

Lastly, by the fact that $x - \sin x \leq \frac{x^3}{6}$ if $x \geq 0$, then we can have

$$\theta(\mathbf{w} + \mathbf{r}, \mathbf{w}) - \sin \theta(\mathbf{w} + \mathbf{r}, \mathbf{w}) = \mathcal{O}\left(\frac{\|\mathbf{r}\|_2^3}{\|\mathbf{w}\|_2^3}\right). \quad (66)$$

Part II: Control the Ratio $\left|1 - \frac{\|\mathbf{w}\|_2}{\|\mathbf{w} + \mathbf{r}\|_2}\right|$

We can compute

$$\left|1 - \frac{\|\mathbf{w}\|_2}{\|\mathbf{w} + \mathbf{r}\|_2}\right| = \left|1 - \frac{1}{\sqrt{1 + 2\alpha + \beta}}\right| = \left|1 - \Theta\left(1 - \alpha - \frac{\beta}{2}\right)\right| = \mathcal{O}\left(\frac{\|\mathbf{r}\|_2}{\|\mathbf{w}\|_2}\right). \quad (67)$$

where the second equality follows from the first order binomial approximation $\frac{1}{\sqrt{1+x}} = \Theta\left(1 - \frac{x}{2}\right)$ if $|x| \ll 1$ and we have $\frac{\|\mathbf{r}\|_2}{\|\mathbf{w}\|_2} = \mathcal{O}\left(\frac{1}{\kappa r^*}\right)$ by Eq. (63). By consequence, we can have

$$\frac{\|\mathbf{w}\|_2}{\|\mathbf{w} + \mathbf{r}\|_2} \leq 1 + \mathcal{O}\left(\frac{\|\mathbf{r}\|_2}{\|\mathbf{w}\|_2}\right). \quad (68)$$

Now, we can upper bound Eq. (62) in terms of Frobenius norm by triangle inequality, i.e.

$$\|\Psi(t)\|_F \leq \|\Psi_1(t)\|_F + \|\Psi_2(t)\|_F + \|\Psi_3(t)\|_F.$$

For $\|\Psi_1(t)\|_F$, we have

$$\begin{aligned}
 \|\Psi_1(t)\|_F &\leq \max_{1 \leq m \leq k} \frac{\|\tilde{\mathbf{w}}_m^\natural\|_2 \sin[\theta(\tilde{\mathbf{w}}_m^\natural + \mathbf{r}_{t,1}, \tilde{\mathbf{w}}_m^\natural)]}{\|\tilde{\mathbf{w}}_m^\natural + \mathbf{r}_{t,m}\|_2} \|\mathbf{A}_t \mathbf{B}_t - \Delta\|_F \\
 &\leq \max_{1 \leq m \leq k} \left(1 + \mathcal{O}\left(\frac{\|\mathbf{r}_{t,m}\|_2}{\|\tilde{\mathbf{w}}_m^\natural\|_2}\right)\right) \mathcal{O}\left(\frac{\|\mathbf{r}_{t,m}\|_2}{\|\tilde{\mathbf{w}}_m^\natural\|_2}\right) \|\mathbf{A}_t \mathbf{B}_t - \Delta\|_F \quad [\text{by Eq. (64) and Eq. (68)}] \\
 &\leq \mathcal{O}\left(\max_{1 \leq m \leq k} \frac{\|\mathbf{r}_{t,m}\|_2}{\|\tilde{\mathbf{w}}_m^\natural\|_2}\right) \|\mathbf{A}_t \mathbf{B}_t - \Delta\|_F \leq \mathcal{O}\left(\frac{1}{\kappa r^*}\right) \|\mathbf{A}_t \mathbf{B}_t - \Delta\|_F. \quad [\text{by Eq. (63)}]
 \end{aligned}$$

For $\|\Psi_2(t)\|_F$, we have

$$\begin{aligned}
 \|\Psi_2(t)\|_F &\leq \|\widetilde{\mathbf{W}}^\natural\|_{op} \|\mathbf{D}_1(t) - \mathbf{D}_2(t)\|_F \\
 &= \|\widetilde{\mathbf{W}}^\natural\|_{op} \sqrt{\sum_{m=1}^k \left(\sin[\theta(\tilde{\mathbf{w}}_m^\natural + \mathbf{r}_{t,m}, \tilde{\mathbf{w}}_m^\natural)] - \theta(\tilde{\mathbf{w}}_m^\natural + \mathbf{r}_{t,m}, \tilde{\mathbf{w}}_m^\natural)\right)^2} \\
 &\leq \|\widetilde{\mathbf{W}}^\natural\|_{op} \sqrt{\sum_{m=1}^k \mathcal{O}\left(\frac{\|\mathbf{r}_{t,m}\|_2^3}{\|\tilde{\mathbf{w}}_m^\natural\|_2^3}\right)} \quad [\text{by Eq. (66)}] \\
 &\leq \|\widetilde{\mathbf{W}}^\natural\|_{op} \sqrt{\sum_{m=1}^k \mathcal{O}\left(\|\mathbf{r}_{t,m}\|_2 \max_{1 \leq i \leq k} \frac{\|\mathbf{r}_{t,i}\|_2^2}{\|\tilde{\mathbf{w}}_i^\natural\|_2^3}\right)} \\
 &= \mathcal{O}\left(\|\widetilde{\mathbf{W}}^\natural\|_{op} \max_{1 \leq i \leq k} \frac{\|\mathbf{r}_{t,i}\|_2^2}{\|\tilde{\mathbf{w}}_i^\natural\|_2^3}\right) \|\mathbf{A}_t \mathbf{B}_t - \Delta\|_F \\
 &\leq \mathcal{O}\left(\frac{\|\widetilde{\mathbf{W}}^\natural\|_{op}}{\min_{1 \leq i \leq k} \|\tilde{\mathbf{w}}_i^\natural\|_2} \max_{1 \leq i \leq k} \frac{\|\mathbf{r}_{t,i}\|_2^2}{\|\tilde{\mathbf{w}}_i^\natural\|_2^3}\right) \|\mathbf{A}_t \mathbf{B}_t - \Delta\|_F \\
 &\leq \mathcal{O}\left(\frac{1}{(\kappa r^*)^2}\right) \|\mathbf{A}_t \mathbf{B}_t - \Delta\|_F. \quad [\text{by Assumption 4.1 and Eq. (63)}]
 \end{aligned}$$

For $\|\Psi_3(t)\|_F$, we have

$$\begin{aligned}
 \|\Psi_3(t)\|_F &\leq \|\widetilde{\mathbf{W}}^\natural\|_{op} \|\mathbf{D}_1(t)(\mathbf{I}_k - \mathbf{D}_3(t))\|_F \\
 &\leq \|\widetilde{\mathbf{W}}^\natural\|_{op} \sqrt{\sum_{m=1}^k \left[\left(1 - \frac{\|\tilde{\mathbf{w}}_m^\natural\|_2}{\|\tilde{\mathbf{w}}_m^\natural + \mathbf{r}_{t,m}\|_2}\right) \sin[\theta(\tilde{\mathbf{w}}_m^\natural + \mathbf{r}_{t,m}, \tilde{\mathbf{w}}_m^\natural)] \right]^2} \\
 &\leq \|\widetilde{\mathbf{W}}^\natural\|_{op} \max_{1 \leq m \leq k} \left| \left(1 - \frac{\|\tilde{\mathbf{w}}_m^\natural\|_2}{\|\tilde{\mathbf{w}}_m^\natural + \mathbf{r}_{t,m}\|_2}\right) \right| \times \sqrt{\sum_{m=1}^k \sin^2[\theta(\tilde{\mathbf{w}}_m^\natural + \mathbf{r}_{t,m}, \tilde{\mathbf{w}}_m^\natural)]} \\
 &\leq \|\widetilde{\mathbf{W}}^\natural\|_{op} \max_{1 \leq m \leq k} \mathcal{O}\left(\frac{\|\mathbf{r}_{t,m}\|_2}{\|\tilde{\mathbf{w}}_m^\natural\|_2}\right) \sqrt{\sum_{m=1}^k \mathcal{O}\left(\frac{\|\mathbf{r}_{t,m}\|_2^2}{\|\tilde{\mathbf{w}}_m^\natural\|_2^2}\right)} \quad [\text{by Eq. (64) and Eq. (67)}] \\
 &\leq \|\widetilde{\mathbf{W}}^\natural\|_{op} \max_{1 \leq m \leq k} \mathcal{O}\left(\frac{\|\mathbf{r}_{t,m}\|_2}{\|\tilde{\mathbf{w}}_m^\natural\|_2}\right) \sqrt{\mathcal{O}\left(\frac{\sum_{m=1}^k \|\mathbf{r}_{t,m}\|_2^2}{\min_{1 \leq i \leq k} \|\tilde{\mathbf{w}}_i^\natural\|_2^2}\right)} \quad [\text{due to the positivity of } \|\mathbf{r}_{t,m}\|_2] \\
 &= \mathcal{O}\left(\frac{\|\widetilde{\mathbf{W}}^\natural\|_{op}}{\min_{1 \leq i \leq k} \|\tilde{\mathbf{w}}_i^\natural\|_2} \max_{1 \leq m \leq k} \frac{\|\mathbf{r}_{t,m}\|_2}{\|\tilde{\mathbf{w}}_m^\natural\|_2}\right) \|\mathbf{A}_t \mathbf{B}_t - \Delta\|_F \\
 &\leq \mathcal{O}\left(\frac{1}{\kappa r^*}\right) \|\mathbf{A}_t \mathbf{B}_t - \Delta\|_F. \quad [\text{by Assumption 4.1 and Eq. (63)}]
 \end{aligned}$$

Combine the above upper bounds together, we can obtain

$$\frac{\|\Psi(t)\|_F}{\|\mathbf{A}_t \mathbf{B}_t - \Delta\|_F} \leq \mathcal{O}\left(\frac{1}{\kappa r^*}\right),$$

which completes the proof. \square

D.1.2. CONCENTRATION OF EMPIRICAL GRADIENTS

In this part, we aim to provide the concentration of empirical gradient $\mathbf{J}_{\mathbf{W}_t} := \mathbf{\Gamma}_{1,t} - \mathbf{\Gamma}_{2,t} \in \mathbb{R}^{d \times k}$ in Frobenius norm. Recall $\mathbf{W}_t := \mathbf{W}^{\natural} + \mathbf{A}_t \mathbf{B}_t$ and $\mathbf{w}_{t,m}$ is the corresponding m -th column of \mathbf{W}_t , denote $\tilde{x}_{i,j}$ as the j -th element of $\tilde{\mathbf{x}}_i$, for notational simplicity, we define each element of $\mathbf{J}_{\mathbf{W}_t} := \mathbf{\Gamma}_{1,t} - \mathbf{\Gamma}_{2,t}$ as

$$c_{t,m}^j(\tilde{\mathbf{x}}_i) := \left(\sigma(\tilde{\mathbf{x}}_i^\top \tilde{\mathbf{w}}_m^{\natural}) - \sigma(\tilde{\mathbf{x}}_i^\top \mathbf{w}_{t,m}^{\natural}) \right) \sigma'(\tilde{\mathbf{x}}_i^\top \mathbf{w}_{t,m}^{\natural}) \tilde{x}_{i,j} \in \mathbb{R}, \quad \text{for } 1 \leq m \leq k, 1 \leq i \leq N, 1 \leq j \leq d,$$

Then, we can write $\mathbf{J}_{\mathbf{W}_t}$ in an element-wise way

$$\mathbf{J}_{\mathbf{W}_t} = \frac{1}{N} \sum_{i=1}^N \begin{bmatrix} c_{t,1}^1(\tilde{\mathbf{x}}_i) & \dots & c_{t,k}^1(\tilde{\mathbf{x}}_i) \\ \vdots & \ddots & \vdots \\ c_{t,1}^d(\tilde{\mathbf{x}}_i) & \dots & c_{t,k}^d(\tilde{\mathbf{x}}_i) \end{bmatrix} \in \mathbb{R}^{d \times k},$$

and

$$\left\| \mathbf{J}_{\mathbf{W}_t} - \mathbb{E}_{\tilde{\mathbf{x}}}[\mathbf{J}_{\mathbf{W}_t}] \right\|_F^2 = \sum_{j=1}^d \sum_{m=1}^k \left(\frac{1}{N} \sum_{i=1}^N c_{t,m}^j(\tilde{\mathbf{x}}_i) - \mathbb{E}_{\tilde{\mathbf{x}}}[c_{t,m}^j(\tilde{\mathbf{x}})] \right)^2.$$

Next, we have the following lemma.

Lemma D.3. *For $1 \leq m \leq k, 1 \leq j \leq d$, under assumptions in Section 2.1 for the nonlinear setting, with probability at least $1 - 2C \exp(-N\epsilon^2)$ for a universal constant $C > 0$ and $\epsilon \in (0, 1)$, we have*

$$\left| \frac{1}{N} \sum_{i=1}^N c_{t,m}^j(\tilde{\mathbf{x}}_i) - \mathbb{E}_{\tilde{\mathbf{x}}}[c_{t,m}^j(\tilde{\mathbf{x}})] \right| \leq C^* K^2 \epsilon \|\tilde{\mathbf{w}}_m^{\natural} - \mathbf{w}_{t,m}^{\natural}\|_2,$$

for some absolute constant $C^* > 0$ and $K = \sqrt{8/3}$.

Proof. Since $\tilde{x}_{i,j} \sim \mathcal{N}(0, 1)$ for any $1 \leq m \leq k$ and $1 \leq j \leq d$, then we have that $K := \|\tilde{x}_{i,j}\|_{\psi_2} = \sqrt{8/3}$. By the Orlicz-based definition of subgaussian norm, the subgaussian norm of random variable is identical to its absolute value. Then, for any $\lambda \in \mathbb{R}$, we have the following moment generating function

$$\begin{aligned} & \mathbb{E} \left[\exp \left(\lambda \left| \left(\sigma(\tilde{\mathbf{x}}_i^\top \tilde{\mathbf{w}}_m^{\natural}) - \sigma(\tilde{\mathbf{x}}_i^\top \mathbf{w}_{t,m}^{\natural}) \right) \sigma'(\tilde{\mathbf{x}}_i^\top \mathbf{w}_{t,m}^{\natural}) \right| \right) \right] \\ & \leq \mathbb{E} \left[\exp \left(\lambda \left| \langle \tilde{\mathbf{x}}_i, \tilde{\mathbf{w}}_m^{\natural} - \mathbf{w}_{t,m}^{\natural} \rangle \right| \right) \right] \quad \text{[by Lipschitz continuity of } \sigma \text{ and } \sigma'] \\ & \leq \mathbb{E} \left[\exp \left((C^*)^2 \lambda^2 \left\| \langle \tilde{\mathbf{x}}_i, \tilde{\mathbf{w}}_m^{\natural} - \mathbf{w}_{t,m}^{\natural} \rangle \right\|_{\psi_2}^2 \right) \right], \quad \text{[by subgaussian property]} \end{aligned}$$

for some constant $C^* > 0$, which implies

$$\left\| \left(\sigma(\tilde{\mathbf{x}}_i^\top \tilde{\mathbf{w}}_m^{\natural}) - \sigma(\tilde{\mathbf{x}}_i^\top \mathbf{w}_{t,m}^{\natural}) \right) \sigma'(\tilde{\mathbf{x}}_i^\top \mathbf{w}_{t,m}^{\natural}) \right\|_{\psi_2}^2 \leq (C^*)^2 \left\| \langle \tilde{\mathbf{x}}_i, \tilde{\mathbf{w}}_m^{\natural} - \mathbf{w}_{t,m}^{\natural} \rangle \right\|_{\psi_2}^2 = (C^* K)^2 \|\tilde{\mathbf{w}}_m^{\natural} - \mathbf{w}_{t,m}^{\natural}\|_2^2,$$

where the last inequality follows from the fact that $\|X\|_{\psi_2} = Ks$ if $X \sim \mathcal{N}(0, s^2)$. Therefore, by Vershynin (2018, Lemma 2.7.7), this implies $c_{t,m}^j(\tilde{\mathbf{x}}_i)$ is sub-exponential with

$$B_{t,m} := \|c_{t,m}^j(\tilde{\mathbf{x}})\|_{\psi_1} \leq \|\tilde{x}_{i,j}\|_{\psi_2} \left\| \left(\sigma(\tilde{\mathbf{x}}_i^\top \tilde{\mathbf{w}}_m^{\natural}) - \sigma(\tilde{\mathbf{x}}_i^\top \mathbf{w}_{t,m}^{\natural}) \right) \sigma'(\tilde{\mathbf{x}}_i^\top \mathbf{w}_{t,m}^{\natural}) \right\|_{\psi_2} \leq C^* K^2 \|\tilde{\mathbf{w}}_m^{\natural} - \mathbf{w}_{t,m}^{\natural}\|_2. \quad (69)$$

Then, let $\epsilon_{t,m} = C^* K^2 \epsilon \|\tilde{\mathbf{w}}_m^{\mathfrak{h}} - \mathbf{w}_{t,m}^{\mathfrak{h}}\|_2$ for $\epsilon \in (0, 1)$, we can apply Bernstein's inequality for sub-exponential variables [Vershynin \(2018, Corollary 2.8.3\)](#)

$$\begin{aligned} & \mathbb{P} \left(\left| \frac{1}{N} \sum_{i=1}^N c_{t,m}^j(\tilde{\mathbf{x}}_i) - \mathbb{E}_{\tilde{\mathbf{x}}} [c_{t,m}^j(\tilde{\mathbf{x}})] \right| \geq \epsilon_{t,m} \right) \\ & \leq 2C \exp \left(-N \min \left\{ \frac{\epsilon_{t,m}}{B_{t,m}}, \frac{\epsilon_{t,m}^2}{B_{t,m}^2} \right\} \right) \quad [\text{for some constant } C > 0] \\ & \leq 2C \exp(-N\epsilon^2) . \quad [\text{by Eq. (69) and } \epsilon \in (0, 1)] \end{aligned}$$

□

Theorem D.4. Suppose $\epsilon \in (0, 1)$, under assumptions in Section 2.1 for the nonlinear setting, then with probability at least $1 - 2Cdk \exp(-N\epsilon^2)$ for a universal constant $C > 0$, we have

$$\left\| \mathbf{J}_{\mathbf{W}_t} - \mathbb{E}_{\tilde{\mathbf{x}}} [\mathbf{J}_{\mathbf{W}_t}] \right\|_{\text{F}} \leq C^* K^2 \sqrt{d\epsilon} \|\mathbf{A}_t \mathbf{B}_t - \Delta\|_{\text{F}} ,$$

for some absolute constant $C^* > 0$ and $K = \sqrt{8/3}$.

Proof. By a union bound argument and Lemma D.3, with probability at least $1 - 2Cdk \exp(-N\epsilon^2)$ for a universal constant $C > 0$, we have

$$\begin{aligned} \left\| \mathbf{J}_{\mathbf{W}_t} - \mathbb{E}_{\tilde{\mathbf{x}}} [\mathbf{J}_{\mathbf{W}_t}] \right\|_{\text{F}}^2 &= \sum_{j=1}^d \sum_{m=1}^k \left(\frac{1}{N} \sum_{i=1}^N c_{t,m}^j(\tilde{\mathbf{x}}_i) - \mathbb{E}_{\tilde{\mathbf{x}}} [c_{t,m}^j(\tilde{\mathbf{x}})] \right)^2 \\ &\leq \sum_{j=1}^d \sum_{m=1}^k \epsilon_{t,m}^2 \\ &\leq \sum_{j=1}^d \sum_{m=1}^k (C^* K^2)^2 \epsilon^2 \|\tilde{\mathbf{w}}_m^{\mathfrak{h}} - \mathbf{w}_{t,m}^{\mathfrak{h}}\|_2^2 \\ &= d(C^* K^2)^2 \epsilon^2 \|\mathbf{A}_t \mathbf{B}_t - \Delta\|_{\text{F}}^2 , \end{aligned}$$

which implies

$$\left\| \mathbf{J}_{\mathbf{W}_t} - \mathbb{E}_{\tilde{\mathbf{x}}} [\mathbf{J}_{\mathbf{W}_t}] \right\|_{\text{F}} \leq C^* K^2 \sqrt{d\epsilon} \|\mathbf{A}_t \mathbf{B}_t - \Delta\|_{\text{F}} ,$$

which finishes the proof. □

Lemma D.5. Recall $\mathbf{G}^{\mathfrak{h}} := -\nabla L(\mathbf{W}^{\mathfrak{h}}) = \mathbf{J}_{\mathbf{W}^{\mathfrak{h}}}$, under Assumption 4.1 and assumptions in Section 2.1 for the nonlinear setting, with (Spectral-init), suppose $\epsilon \leq \frac{\rho}{3C^* K^2 \gamma \sqrt{2dr^*} \kappa}$ for some positive constant $\rho > 0$ and we set $\gamma = 2$, then with probability at least $1 - 2Cdk \exp(-N\epsilon^2)$ for a universal constant $C > 0$, it holds that

$$\|\mathbf{A}_0 \mathbf{B}_0 - \Delta\|_{\text{F}} \leq \rho \lambda_{r^*}^* .$$

Proof. We start with decompose $\|\mathbf{A}_0 \mathbf{B}_0 - \Delta\|_{\text{op}}$ into three components, i.e.

$$\|\mathbf{A}_0 \mathbf{B}_0 - \Delta\|_{\text{op}} \leq \underbrace{\|\mathbf{A}_0 \mathbf{B}_0 - \gamma \mathbf{G}^{\mathfrak{h}}\|_{\text{op}}}_{\text{low-rank approximation error}} + \underbrace{\gamma \|\mathbf{G}^{\mathfrak{h}} - \mathbb{E}_{\tilde{\mathbf{x}}} [\mathbf{G}^{\mathfrak{h}}]\|_{\text{op}}}_{\text{concentration error}} + \underbrace{\|\gamma \mathbb{E}_{\tilde{\mathbf{x}}} [\mathbf{G}^{\mathfrak{h}}] - \Delta\|_{\text{op}}}_{\text{population error}} . \quad (70)$$

First, for the population error, we can use similar technique from Lemma D.2, since $\frac{\|\Delta_m\|_2}{\|\tilde{\mathbf{w}}_m^{\mathfrak{h}}\|_2} = \mathcal{O}(\frac{1}{\kappa r^*})$ by Assumption 4.1 for $1 \leq m \leq k$, we can obtain

$$\frac{\|\mathbb{E}_{\tilde{\mathbf{x}}} [\mathbf{G}^{\mathfrak{h}}] - \frac{1}{2} \Delta\|_{\text{F}}}{\|\Delta\|_{\text{F}}} \leq \mathcal{O} \left(\frac{1}{\kappa r^*} \right) .$$

Next, for the concentration error, following Theorem D.4, we replace \mathbf{W}_t with \mathbf{W}^\natural and then obtain the following concentration with the probability at least $1 - 2Cdk \exp(-N\epsilon^2)$ for a universal constant $C > 0$, we have

$$\|\mathbf{G}^\natural - \mathbb{E}_{\tilde{\mathbf{x}}} [\mathbf{G}^\natural]\|_F \leq \frac{\rho \|\Delta\|_F}{3\sqrt{2}r^*\gamma\kappa} \leq \frac{\rho\sqrt{r^*}\|\Delta\|_{op}}{3\sqrt{2}r^*\gamma\kappa} = \frac{\rho\lambda_{r^*}^*}{3\sqrt{2}r^*\gamma}, \quad (71)$$

where $\epsilon \leq \frac{\rho}{2C^*K^2\gamma\sqrt{2dr^*}\kappa}$ for $\rho > 0$.

Lastly, we can upper bound the $(r^* + 1)$ -th singular value of \mathbf{G}^\natural (with scale parameter γ) which acts as the low-rank approximation error. Due to the randomness contained in \mathbf{G}^\natural , we decompose the $(r^* + 1)$ -th singular value into two components, i.e.

$$\gamma\lambda_{r^*+1}(\mathbf{G}^\natural) \leq \underbrace{|\gamma\lambda_{r^*+1}(\mathbf{G}^\natural) - \lambda_{r^*+1}(\gamma\mathbb{E}_{\tilde{\mathbf{x}}}[\mathbf{G}^\natural])|}_{\text{concentration error}} + \underbrace{\lambda_{r^*+1}(\gamma\mathbb{E}_{\tilde{\mathbf{x}}}[\mathbf{G}^\natural])}_{\text{population error}}.$$

First, for the concentration error, we can obtain

$$\begin{aligned} & |\gamma\lambda_{r^*+1}(\mathbf{G}^\natural) - \lambda_{r^*+1}(\gamma\mathbb{E}_{\tilde{\mathbf{x}}}[\mathbf{G}^\natural])| \\ & \leq \gamma \|\mathbf{G}^\natural - \mathbb{E}_{\tilde{\mathbf{x}}}[\mathbf{G}^\natural]\|_{op} \quad [\text{by Weyl's inequality}] \\ & \leq \gamma \|\mathbf{G}^\natural - \mathbb{E}_{\tilde{\mathbf{x}}}[\mathbf{G}^\natural]\|_F \\ & \leq \frac{\rho\lambda_{r^*}^*}{3\sqrt{2}r^*}. \quad [\text{by Eq. (71)}] \end{aligned}$$

Second, we can obtain the population error as

$$\begin{aligned} \lambda_{r^*+1}(\mathbb{E}_{\tilde{\mathbf{x}}}[\mathbf{G}^\natural]) &= \left| \lambda_{r^*+1}(\mathbb{E}_{\tilde{\mathbf{x}}}[\mathbf{G}^\natural]) - \frac{1}{2}\lambda_{r^*+1}(\Delta) \right| \quad [\text{since Rank}(\Delta) = r^*] \\ &\leq \left\| \mathbb{E}_{\tilde{\mathbf{x}}}[\mathbf{G}^\natural] - \frac{1}{2}\Delta \right\|_{op} \quad [\text{by Weyl's inequality}] \\ &\leq \left\| \mathbb{E}_{\tilde{\mathbf{x}}}[\mathbf{G}^\natural] - \frac{1}{2}\Delta \right\|_F \\ &\leq \mathcal{O}\left(\frac{1}{\kappa r^*}\right) \|\Delta\|_F. \end{aligned}$$

Now we can have

$$\gamma\lambda_{r^*+1}(\mathbf{G}^\natural) \leq \frac{\rho\lambda_{r^*}^*}{3\sqrt{2}r^*} + \mathcal{O}\left(\frac{\|\Delta\|_F}{\kappa r^*}\right) \leq \mathcal{O}\left(\frac{1}{\sqrt{r^*}}\right) \rho\lambda_{r^*}^*. \quad (72)$$

Therefore, combine everything together, recall Eq. (70), we can obtain

$$\begin{aligned} \|\mathbf{A}_0\mathbf{B}_0 - \Delta\|_{op} &\leq \gamma\lambda_{r^*+1}(\mathbf{G}^\natural) + \gamma \|\mathbf{G}^\natural - \mathbb{E}_{\tilde{\mathbf{x}}}[\mathbf{G}^\natural]\|_F + \gamma \left\| \mathbb{E}_{\tilde{\mathbf{x}}}[\mathbf{G}^\natural] - \frac{1}{\gamma}\Delta \right\|_F \\ &\leq \mathcal{O}\left(\frac{1}{\sqrt{r^*}}\right) \rho\lambda_{r^*}^* + \frac{\rho\lambda_{r^*}^*}{3\sqrt{2}r^*} + \mathcal{O}\left(\frac{1}{\sqrt{r^*}}\right) \rho\lambda_{r^*}^* \leq \frac{\rho\lambda_{r^*}^*}{\sqrt{2}r^*}. \end{aligned} \quad (73)$$

Since we work in the exact-rank case $\text{Rank}(\mathbf{A}_t\mathbf{B}_t) \leq r = r^*$ with $\text{Rank}(\Delta) = r^*$, then $\text{Rank}(\mathbf{A}_0\mathbf{B}_0 - \Delta) \leq 2r^*$, this can imply

$$\|\mathbf{A}_0\mathbf{B}_0 - \Delta\|_F \leq \sqrt{2r^*} \|\mathbf{A}_0\mathbf{B}_0 - \Delta\|_{op} \leq \rho\lambda_{r^*}^*,$$

which completes the proof. \square

D.2. Preconditioned Gradient Descent under Spectral Initialization

Recall the loss of LoRA fine-tuning:

$$\tilde{L}(\mathbf{A}_t, \mathbf{B}_t) = \frac{1}{2N} \left\| \sigma \left(\tilde{\mathbf{X}} (\mathbf{W}^\natural + \mathbf{A}_t \mathbf{B}_t) \right) - \sigma \left(\tilde{\mathbf{X}} \tilde{\mathbf{W}}^\natural \right) \right\|_F^2.$$

Then, we employ the following preconditioned gradient updates for LoRA fine-tuning

$$\mathbf{A}_{t+1} = \mathbf{A}_t + \eta \mathbf{J}_{\mathbf{W}_t} \mathbf{B}_t^\top (\mathbf{B}_t \mathbf{B}_t^\top)^{-1}, \quad (74)$$

and

$$\mathbf{B}_{t+1} = \mathbf{B}_t + \eta (\mathbf{A}_t^\top \mathbf{A}_t)^{-1} \mathbf{A}_t^\top \mathbf{J}_{\mathbf{W}_t}. \quad (75)$$

Similar to the linear case, we define the following notations

- SVD of product matrix $\mathbf{A}_t \mathbf{B}_t := \mathcal{U}_t \mathcal{S}_t \mathcal{V}_t^\top$, where $\mathcal{U}_t \in \mathbb{R}^{d \times r}$, $\mathcal{S}_t \in \mathbb{R}^{r^* \times r}$, and $\mathcal{V}_t \in \mathbb{R}^{k \times r}$.
- The left singular matrix of \mathbf{A}_t as $\mathbf{U}_{\mathbf{A}_t} \in \mathbb{R}^{d \times r}$.
- The right singular matrix of \mathbf{B}_t as $\mathbf{V}_{\mathbf{B}_t} \in \mathbb{R}^{k \times r}$.

Lemma D.6. *Under assumptions in Section 2.1 for the nonlinear setting, we update \mathbf{A}_t and \mathbf{B}_t via Eq. (74) and Eq. (75) under spectral initialization (Spectral-init), then we have the following recursion*

$$\begin{aligned} \mathbf{A}_{t+1} \mathbf{B}_{t+1} - \Delta &= (1 - \eta) \mathbf{U}_{\mathbf{A}_t} \mathbf{U}_{\mathbf{A}_t}^\top (\mathbf{A}_t \mathbf{B}_t - \Delta) \mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^\top \\ &\quad + (1 - \eta/2) (\mathbf{I}_d - \mathbf{U}_{\mathbf{A}_t} \mathbf{U}_{\mathbf{A}_t}^\top) (\mathbf{A}_t \mathbf{B}_t - \Delta) \mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^\top \\ &\quad + (1 - \eta/2) \mathbf{U}_{\mathbf{A}_t} \mathbf{U}_{\mathbf{A}_t}^\top (\mathbf{A}_t \mathbf{B}_t - \Delta) (\mathbf{I}_k - \mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^\top) \\ &\quad + (\mathbf{I}_d - \mathbf{U}_{\mathbf{A}_t} \mathbf{U}_{\mathbf{A}_t}^\top) (\mathbf{A}_t \mathbf{B}_t - \Delta) (\mathbf{I}_k - \mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^\top) \\ &\quad + \eta \Xi_t \mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^\top + \eta \mathbf{U}_{\mathbf{A}_t} \mathbf{U}_{\mathbf{A}_t}^\top \Xi_t + \eta^2 \mathbf{J}_{\mathbf{W}_t} \mathcal{V}_t \mathcal{S}_t^{-1} \mathcal{U}_t^\top \mathbf{J}_{\mathbf{W}_t}, \end{aligned} \quad (76)$$

and

$$\Xi_t := \mathbf{J}_{\mathbf{W}_t} - \frac{1}{2} (\mathbf{A}_t \mathbf{B}_t - \Delta).$$

Then, by choosing $\eta \in (0, 1)$, we have the associated upper bound in Frobenius norm

$$\begin{aligned} &\|\mathbf{A}_{t+1} \mathbf{B}_{t+1} - \Delta\|_F \\ &\leq (1 - \eta) \|\mathbf{U}_{\mathbf{A}_t} \mathbf{U}_{\mathbf{A}_t}^\top (\mathbf{A}_t \mathbf{B}_t - \Delta) \mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^\top\|_F \end{aligned} \quad (77)$$

$$\begin{aligned} &\quad + (1 - \eta/2) \|(\mathbf{I}_d - \mathbf{U}_{\mathbf{A}_t} \mathbf{U}_{\mathbf{A}_t}^\top) (\mathbf{A}_t \mathbf{B}_t - \Delta) \mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^\top + \mathbf{U}_{\mathbf{A}_t} \mathbf{U}_{\mathbf{A}_t}^\top (\mathbf{A}_t \mathbf{B}_t - \Delta) (\mathbf{I}_k - \mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^\top)\|_F \\ &\quad + \|(\mathbf{I}_d - \mathbf{U}_{\mathbf{A}_t} \mathbf{U}_{\mathbf{A}_t}^\top) (\mathbf{A}_t \mathbf{B}_t - \Delta) (\mathbf{I}_k - \mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^\top)\|_F \\ &\quad + 2\eta \|\Xi_t\|_F + \eta^2 \|\mathbf{J}_{\mathbf{W}_t} \mathcal{V}_t \mathcal{S}_t^{-1} \mathcal{U}_t^\top \mathbf{J}_{\mathbf{W}_t}\|_F. \end{aligned} \quad (78)$$

Proof. By the preconditioned update in Eq. (74) and Eq. (75), we can construct

$$\begin{aligned}
 \mathbf{A}_{t+1}\mathbf{B}_{t+1} - \Delta &= \mathbf{A}_t\mathbf{B}_t - \Delta \\
 &\quad - \eta \mathbf{J}_{\mathbf{W}_t} \mathbf{B}_t^\top (\mathbf{B}_t \mathbf{B}_t^\top)^{-1} \mathbf{B}_t - \eta \mathbf{A}_t (\mathbf{A}_t^\top \mathbf{A}_t)^{-1} \mathbf{A}_t^\top \mathbf{J}_{\mathbf{W}_t} \\
 &\quad + \eta^2 \mathbf{J}_{\mathbf{W}_t} \mathbf{B}_t^\top (\mathbf{B}_t \mathbf{B}_t^\top)^{-1} (\mathbf{A}_t^\top \mathbf{A}_t)^{-1} \mathbf{A}_t^\top \mathbf{J}_{\mathbf{W}_t} \\
 &= \mathbf{A}_t\mathbf{B}_t - \Delta \\
 &\quad - \eta/2 (\mathbf{A}_t\mathbf{B}_t - \Delta) \mathbf{B}_t^\top (\mathbf{B}_t \mathbf{B}_t^\top)^{-1} \mathbf{B}_t + \eta \Xi_t \mathbf{B}_t^\top (\mathbf{B}_t \mathbf{B}_t^\top)^{-1} \mathbf{B}_t \\
 &\quad - \eta/2 \mathbf{A}_t (\mathbf{A}_t^\top \mathbf{A}_t)^{-1} \mathbf{A}_t^\top (\mathbf{A}_t\mathbf{B}_t - \Delta) + \eta \mathbf{A}_t (\mathbf{A}_t^\top \mathbf{A}_t)^{-1} \mathbf{A}_t^\top \Xi_t \\
 &\quad + \eta^2 \mathbf{J}_{\mathbf{W}_t} \mathbf{B}_t^\top (\mathbf{B}_t \mathbf{B}_t^\top)^{-1} (\mathbf{A}_t^\top \mathbf{A}_t)^{-1} \mathbf{A}_t^\top \mathbf{J}_{\mathbf{W}_t} \\
 &= \mathbf{A}_t\mathbf{B}_t - \Delta \\
 &\quad - \eta/2 (\mathbf{A}_t\mathbf{B}_t - \Delta) \mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^\top + \eta \Xi_t \mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^\top \\
 &\quad - \eta/2 \mathbf{U}_{\mathbf{A}_t} \mathbf{U}_{\mathbf{A}_t}^\top (\mathbf{A}_t\mathbf{B}_t - \Delta) + \eta \mathbf{U}_{\mathbf{A}_t} \mathbf{U}_{\mathbf{A}_t}^\top \Xi_t \\
 &\quad + \eta^2 \mathbf{J}_{\mathbf{W}_t} \mathcal{V}_t \mathcal{S}_t^{-1} \mathcal{U}_t^\top \mathbf{J}_{\mathbf{W}_t}, \quad \text{[by pseudo inverse theorem and Jia et al. (2024, Lemma 14)]}
 \end{aligned}$$

from our choice on $\Xi_t = \mathbf{J}_{\mathbf{W}_t} - \frac{1}{2} (\mathbf{A}_t\mathbf{B}_t - \Delta)$. We can continue to expand

$$\begin{aligned}
 \mathbf{A}_{t+1}\mathbf{B}_{t+1} - \Delta &= (\mathbf{I}_d - \mathbf{U}_{\mathbf{A}_t} \mathbf{U}_{\mathbf{A}_t}^\top + \mathbf{U}_{\mathbf{A}_t} \mathbf{U}_{\mathbf{A}_t}^\top) (\mathbf{A}_t\mathbf{B}_t - \Delta) (\mathbf{I}_d - \mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^\top + \mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^\top) \\
 &\quad - \eta/2 (\mathbf{I}_d - \mathbf{U}_{\mathbf{A}_t} \mathbf{U}_{\mathbf{A}_t}^\top + \mathbf{U}_{\mathbf{A}_t} \mathbf{U}_{\mathbf{A}_t}^\top) (\mathbf{A}_t\mathbf{B}_t - \Delta) \mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^\top \\
 &\quad - \eta/2 \mathbf{U}_{\mathbf{A}_t} \mathbf{U}_{\mathbf{A}_t}^\top (\mathbf{A}_t\mathbf{B}_t - \Delta) (\mathbf{I}_d - \mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^\top + \mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^\top) \\
 &\quad + \eta \Xi_t \mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^\top + \eta \mathbf{U}_{\mathbf{A}_t} \mathbf{U}_{\mathbf{A}_t}^\top \Xi_t + \eta^2 \mathbf{J}_{\mathbf{W}_t} \mathcal{V}_t \mathcal{S}_t^{-1} \mathcal{U}_t^\top \mathbf{J}_{\mathbf{W}_t} \\
 &= (1 - \eta) \mathbf{U}_{\mathbf{A}_t} \mathbf{U}_{\mathbf{A}_t}^\top (\mathbf{A}_t\mathbf{B}_t - \Delta) \mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^\top \\
 &\quad + (1 - \eta/2) (\mathbf{I}_d - \mathbf{U}_{\mathbf{A}_t} \mathbf{U}_{\mathbf{A}_t}^\top) (\mathbf{A}_t\mathbf{B}_t - \Delta) \mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^\top \\
 &\quad + (1 - \eta/2) \mathbf{U}_{\mathbf{A}_t} \mathbf{U}_{\mathbf{A}_t}^\top (\mathbf{A}_t\mathbf{B}_t - \Delta) (\mathbf{I}_d - \mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^\top) \\
 &\quad + (\mathbf{I}_d - \mathbf{U}_{\mathbf{A}_t} \mathbf{U}_{\mathbf{A}_t}^\top) (\mathbf{A}_t\mathbf{B}_t - \Delta) (\mathbf{I}_d - \mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^\top) \\
 &\quad + \eta \Xi_t \mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^\top + \eta \mathbf{U}_{\mathbf{A}_t} \mathbf{U}_{\mathbf{A}_t}^\top \Xi_t + \eta^2 \mathbf{J}_{\mathbf{W}_t} \mathcal{V}_t \mathcal{S}_t^{-1} \mathcal{U}_t^\top \mathbf{J}_{\mathbf{W}_t}.
 \end{aligned}$$

Based on the above formulation, suppose $\eta \in (0, 1)$, we can derive the following upper bound by triangle inequality

$$\begin{aligned}
 &\|\mathbf{A}_{t+1}\mathbf{B}_{t+1} - \Delta\|_F \\
 &\leq \|(1 - \eta) \mathbf{U}_{\mathbf{A}_t} \mathbf{U}_{\mathbf{A}_t}^\top (\mathbf{A}_t\mathbf{B}_t - \Delta) \mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^\top\|_F \\
 &\quad + \|(1 - \eta/2) (\mathbf{I}_d - \mathbf{U}_{\mathbf{A}_t} \mathbf{U}_{\mathbf{A}_t}^\top) (\mathbf{A}_t\mathbf{B}_t - \Delta) \mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^\top\|_F \\
 &\quad + \|(1 - \eta/2) \mathbf{U}_{\mathbf{A}_t} \mathbf{U}_{\mathbf{A}_t}^\top (\mathbf{A}_t\mathbf{B}_t - \Delta) (\mathbf{I}_d - \mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^\top)\|_F \\
 &\quad + \|(\mathbf{I}_d - \mathbf{U}_{\mathbf{A}_t} \mathbf{U}_{\mathbf{A}_t}^\top) (\mathbf{A}_t\mathbf{B}_t - \Delta) (\mathbf{I}_d - \mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^\top)\|_F \\
 &\quad + \eta \|\Xi_t \mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^\top\|_F + \eta \|\mathbf{U}_{\mathbf{A}_t} \mathbf{U}_{\mathbf{A}_t}^\top \Xi_t\|_F + \eta^2 \|\mathbf{J}_{\mathbf{W}_t} \mathcal{V}_t \mathcal{S}_t^{-1} \mathcal{U}_t^\top \mathbf{J}_{\mathbf{W}_t}\|_F \quad \text{[by triangle inequality]} \\
 &\leq (1 - \eta) \|\mathbf{U}_{\mathbf{A}_t} \mathbf{U}_{\mathbf{A}_t}^\top (\mathbf{A}_t\mathbf{B}_t - \Delta) \mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^\top\|_F \\
 &\quad + (1 - \eta/2) \|(\mathbf{I}_d - \mathbf{U}_{\mathbf{A}_t} \mathbf{U}_{\mathbf{A}_t}^\top) (\mathbf{A}_t\mathbf{B}_t - \Delta) \mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^\top + \mathbf{U}_{\mathbf{A}_t} \mathbf{U}_{\mathbf{A}_t}^\top (\mathbf{A}_t\mathbf{B}_t - \Delta) (\mathbf{I}_d - \mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^\top)\|_F \\
 &\quad + \|(\mathbf{I}_d - \mathbf{U}_{\mathbf{A}_t} \mathbf{U}_{\mathbf{A}_t}^\top) (\mathbf{A}_t\mathbf{B}_t - \Delta) (\mathbf{I}_d - \mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^\top)\|_F \\
 &\quad + 2\eta \|\Xi_t\|_F + \eta^2 \|\mathbf{J}_{\mathbf{W}_t} \mathcal{V}_t \mathcal{S}_t^{-1} \mathcal{U}_t^\top \mathbf{J}_{\mathbf{W}_t}\|_F, \quad \text{[since } \eta \in (0, 1)\text{]}
 \end{aligned}$$

which proves the claim. \square

In order to derive the convergence rate of $\|\mathbf{A}_{t+1}\mathbf{B}_{t+1} - \Delta\|_F$ in the above terms, we need to provide the estimation of the

following four terms

$$\begin{aligned} & \|\Xi_t\|_F, \quad \|\mathbf{J}_{\mathbf{W}_t} \mathcal{V}_t \mathcal{S}_t^{-1} \mathcal{U}_t^\top \mathbf{J}_{\mathbf{W}_t}\|_F, \\ & \left\| (\mathbf{I}_d - \mathbf{U}_{\mathbf{A}_t} \mathbf{U}_{\mathbf{A}_t}^\top) (\mathbf{A}_t \mathbf{B}_t - \Delta) \mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^\top + \mathbf{U}_{\mathbf{A}_t} \mathbf{U}_{\mathbf{A}_t}^\top (\mathbf{A}_t \mathbf{B}_t - \Delta) (\mathbf{I}_k - \mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^\top) \right\|_F, \\ & \left\| (\mathbf{I}_d - \mathbf{U}_{\mathbf{A}_t} \mathbf{U}_{\mathbf{A}_t}^\top) (\mathbf{A}_t \mathbf{B}_t - \Delta) (\mathbf{I}_k - \mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^\top) \right\|_F. \end{aligned}$$

which are important elements in Eq. (77). We firstly prove the upper bound for $\|\Xi_t\|_F$ and $\|\mathbf{J}_{\mathbf{W}_t} \mathcal{V}_t \mathcal{S}_t^{-1} \mathcal{U}_t^\top \mathbf{J}_{\mathbf{W}_t}\|_F$ since they are relatively straightforward. After that, we will handle with the remaining three terms which are the most technical part. All of these three terms rely on the condition $\|\mathbf{A}_t \mathbf{B}_t - \Delta\|_F \leq \rho \lambda_{r^*}^*$ and we will prove it by induction finally in Theorem D.10.

Lemma D.7. *For a positive constant $\rho \in (0, 1)$, suppose $\epsilon \leq \frac{\rho}{3C^* K^2 \gamma \sqrt{2d} r^* \kappa}$ with $\gamma = 2$, assume $\|\mathbf{A}_t \mathbf{B}_t - \Delta\|_F \leq \rho \lambda_{r^*}^*$, under assumptions in Section 2.1 for the nonlinear setting and Assumption 4.1, then with probability at least $1 - 2Cdk \exp(-\epsilon^2 N)$ for a universal constant $C > 0$, we have*

$$\|\Xi_t\|_F \leq \left(\mathcal{O}\left(\frac{1}{\kappa r^*}\right) + C^* K^2 \sqrt{d} \epsilon \right) \|\mathbf{A}_t \mathbf{B}_t - \Delta\|_F,$$

Proof. Recall $\Xi_t := \mathbf{J}_{\mathbf{W}_t} - \frac{1}{2} (\mathbf{A}_t \mathbf{B}_t - \Delta)$ from Lemma D.6, then with probability at least $1 - 2Cdk \exp(-\epsilon^2 N)$ for a universal constant $C > 0$, we have

$$\begin{aligned} \|\Xi_t\|_F &= \left\| \frac{1}{2} (\mathbf{A}_t \mathbf{B}_t - \Delta) - \mathbb{E}_{\tilde{\mathbf{x}}} [\mathbf{J}_{\mathbf{W}_t}] + \mathbb{E}_{\tilde{\mathbf{x}}} [\mathbf{J}_{\mathbf{W}_t}] - \mathbf{J}_{\mathbf{W}_t} \right\|_F \\ &\leq \underbrace{\left\| \frac{1}{2} (\mathbf{A}_t \mathbf{B}_t - \Delta) - \mathbb{E}_{\tilde{\mathbf{x}}} [\mathbf{J}_{\mathbf{W}_t}] \right\|_F}_{\text{population error}} + \underbrace{\|\mathbb{E}_{\tilde{\mathbf{x}}} [\mathbf{J}_{\mathbf{W}_t}] - \mathbf{J}_{\mathbf{W}_t}\|_F}_{\text{concentration error}} \\ &= \frac{\|\Psi(t)\|_F}{\|\mathbf{A}_t \mathbf{B}_t - \Delta\|_F} \|\mathbf{A}_t \mathbf{B}_t - \Delta\|_F + \|\mathbb{E}_{\tilde{\mathbf{x}}} [\mathbf{J}_{\mathbf{W}_t}] - \mathbf{J}_{\mathbf{W}_t}\|_F \quad [\text{by Lemma D.2}] \\ &\leq \left(\mathcal{O}\left(\frac{1}{\kappa r^*}\right) + C^* K^2 \sqrt{d} \epsilon \right) \|\mathbf{A}_t \mathbf{B}_t - \Delta\|_F, \quad [\text{by Lemma D.2 and Theorem D.4}] \end{aligned}$$

which completes the proof. \square

Lemma D.8. *Under assumptions in Section 2.1 for the nonlinear setting, suppose $\|\mathbf{A}_t \mathbf{B}_t - \Delta\|_F \leq \rho \lambda_{r^*}^*$ for a positive constant $\rho > 0$, with probability at least $1 - 2C \exp(-\epsilon^2 N)$ for some constants $C > 0$, it holds that*

$$\|\mathbf{J}_{\mathbf{W}_t} \mathcal{V}_t \mathcal{S}_t^{-1} \mathcal{U}_t^\top \mathbf{J}_{\mathbf{W}_t}\|_F \leq (1 + \epsilon)^2 \frac{\rho}{1 - \rho} \|\mathbf{A}_t \mathbf{B}_t - \Delta\|_F.$$

Proof. First, with probability at least $1 - 2C \exp(-\epsilon^2 N)$ for some constants $C > 0$, we can derive

$$\begin{aligned} \|\mathbf{J}_{\mathbf{W}_t} \mathcal{V}_t \mathcal{S}_t^{-1} \mathcal{U}_t^\top \mathbf{J}_{\mathbf{W}_t}\|_F &\leq \|\mathbf{J}_{\mathbf{W}_t}\|_F^2 \|\mathcal{V}_t \mathcal{S}_t^{-1} \mathcal{U}_t^\top\|_{op} \\ &= \frac{\left\| \frac{1}{N} \widetilde{\mathbf{X}}^\top \left(\sigma \left(\widetilde{\mathbf{X}} (\mathbf{W}^\natural + \mathbf{A}_t \mathbf{B}_t) \right) - \sigma \left(\widetilde{\mathbf{X}} \widetilde{\mathbf{W}}^\natural \right) \right) \odot \sigma' \left(\widetilde{\mathbf{X}} (\mathbf{W}^\natural + \mathbf{A}_t \mathbf{B}_t) \right) \right\|_F^2}{\lambda_r(\mathbf{A}_t \mathbf{B}_t)} \\ &\leq \frac{\left\| \frac{1}{N} \widetilde{\mathbf{X}}^\top \widetilde{\mathbf{X}} (\mathbf{A}_t \mathbf{B}_t - \Delta) \right\|_F^2}{\lambda_r(\mathbf{A}_t \mathbf{B}_t)} \quad [\text{by Lipschitz continuity of } \sigma, \sigma'] \\ &\leq \left(\frac{1}{N} \lambda_1^2(\widetilde{\mathbf{X}}) \right)^2 \frac{\|\mathbf{A}_t \mathbf{B}_t - \Delta\|_F^2}{\lambda_r(\mathbf{A}_t \mathbf{B}_t)} \\ &\leq (1 + \epsilon)^2 \frac{\rho}{1 - \rho} \|\mathbf{A}_t \mathbf{B}_t - \Delta\|_F, \quad [\text{by concentration of operator norm}] \end{aligned}$$

where the last equality follows from $r = r^*$ and

$$\lambda_r(\mathbf{A}_t \mathbf{B}_t) \geq \lambda_{r^*}(\Delta) - \|\mathbf{A}_t \mathbf{B}_t - \Delta\|_F \geq (1 - \rho) \lambda_{r^*}(\Delta).$$

□

With Lemma D.7 and Lemma D.8, now we can prove for the other three terms.

Lemma D.9. Suppose $\|\mathbf{A}_t \mathbf{B}_t - \Delta\|_F \leq \rho \lambda_{r^*}$ with a positive constant $\rho \in [0, 1/4]$, then it holds that

$$\|(\mathbf{I}_d - \mathbf{U}_{\mathbf{A}_t} \mathbf{U}_{\mathbf{A}_t}^\top) \Delta (\mathbf{I}_k - \mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^\top)\|_F \leq \frac{\rho}{\sqrt{1 - 8\rho^2}} \|\mathbf{A}_t \mathbf{B}_t - \Delta\|_F,$$

and

$$\|(\mathbf{I}_d - \mathbf{U}_{\mathbf{A}_t} \mathbf{U}_{\mathbf{A}_t}^\top) \Delta \mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^\top + \mathbf{U}_{\mathbf{A}_t} \mathbf{U}_{\mathbf{A}_t}^\top \Delta (\mathbf{I}_k - \mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^\top)\|_F \leq \|\mathbf{A}_t \mathbf{B}_t - \Delta\|_F.$$

Proof. First, we recall

$$\mathbf{Z}_t = \begin{bmatrix} \mathbf{A}_t \\ \mathbf{B}_t^\top \end{bmatrix}, \quad \underline{\mathbf{Z}}_t = \begin{bmatrix} \mathbf{A}_t \\ -\mathbf{B}_t^\top \end{bmatrix},$$

and define a preconditioned operator \mathcal{P} and symmetrized downstream feature shift matrix $\hat{\Delta}$ as

$$\mathcal{P}(\mathbf{Z}_t) := \begin{bmatrix} \mathbf{A}_t (\mathbf{A}_t^\top \mathbf{A}_t)^{-1} \\ \mathbf{B}_t^\top (\mathbf{B}_t \mathbf{B}_t^\top)^{-1} \end{bmatrix}, \quad \mathcal{P}(\underline{\mathbf{Z}}_t) := \begin{bmatrix} \mathbf{A}_t (\mathbf{A}_t^\top \mathbf{A}_t)^{-1} \\ -\mathbf{B}_t^\top (\mathbf{B}_t \mathbf{B}_t^\top)^{-1} \end{bmatrix}, \quad \hat{\Delta} := \begin{bmatrix} \mathbf{0}_{d \times d} & \Delta \\ \Delta^\top & \mathbf{0}_{k \times k} \end{bmatrix}.$$

Next, we observe that

$$\frac{1}{2} (\mathbf{Z}_t \mathbf{Z}_t^\top - \underline{\mathbf{Z}}_t \underline{\mathbf{Z}}_t^\top) - \hat{\Delta} = \begin{bmatrix} \mathbf{0}_{d \times d} & \mathbf{A}_t \mathbf{B}_t - \Delta \\ (\mathbf{A}_t \mathbf{B}_t - \Delta)^\top & \mathbf{0}_{k \times k} \end{bmatrix},$$

leading to

$$\left\| \frac{1}{2} (\mathbf{Z}_t \mathbf{Z}_t^\top - \underline{\mathbf{Z}}_t \underline{\mathbf{Z}}_t^\top) - \hat{\Delta} \right\|_{op} = \|\mathbf{A}_t \mathbf{B}_t - \Delta\|_{op}, \quad \left\| \frac{1}{2} (\mathbf{Z}_t \mathbf{Z}_t^\top - \underline{\mathbf{Z}}_t \underline{\mathbf{Z}}_t^\top) - \hat{\Delta} \right\|_F = \sqrt{2} \|\mathbf{A}_t \mathbf{B}_t - \Delta\|_F.$$

Based on the compact SVD of Δ in Eq. (1), we can write out the eigendecomposition of $\hat{\Delta}$ as

$$\hat{\Delta} = \begin{bmatrix} \Phi & \underline{\Phi} \end{bmatrix} \begin{bmatrix} \mathbf{S}^* & \mathbf{0}_{r^* \times r^*} \\ \mathbf{0}_{r^* \times r^*} & -\mathbf{S}^* \end{bmatrix} \begin{bmatrix} \Phi & \underline{\Phi} \end{bmatrix}^\top = \Phi \mathbf{S}^* \Phi^\top - \underline{\Phi} \mathbf{S}^* \underline{\Phi}^\top, \quad \text{where } \Phi = \frac{1}{\sqrt{2}} \begin{bmatrix} \mathbf{U} \\ \mathbf{V} \end{bmatrix}, \underline{\Phi} = \frac{1}{\sqrt{2}} \begin{bmatrix} \mathbf{U} \\ -\mathbf{V} \end{bmatrix}. \quad (80)$$

Notice that we can also obtain the SVD of $\hat{\Delta}$ as

$$\hat{\Delta} = \hat{\mathbf{U}} \hat{\mathbf{S}} \hat{\mathbf{V}}^\top = \underbrace{\begin{bmatrix} \Phi & \underline{\Phi} \end{bmatrix}}_{:= \hat{\mathbf{U}}} \underbrace{\begin{bmatrix} \mathbf{S}^* & \mathbf{0}_{r^* \times r^*} \\ \mathbf{0}_{r^* \times r^*} & -\mathbf{S}^* \end{bmatrix}}_{\hat{\mathbf{S}}} \underbrace{\begin{bmatrix} \Phi & -\underline{\Phi} \end{bmatrix}^\top}_{:= \hat{\mathbf{V}}^\top}. \quad (81)$$

Notice that $\hat{\Delta}$ is a low-rank matrix with $\text{rank}(2r^*)$ because of $\text{Rank}(\Delta) = r^*$. If $\frac{1}{2} (\mathbf{Z}_t \mathbf{Z}_t^\top - \underline{\mathbf{Z}}_t \underline{\mathbf{Z}}_t^\top)$ recovers $\hat{\Delta}$, this indicates that the top- $2r^*$ subspace of $\frac{1}{2} (\mathbf{Z}_t \mathbf{Z}_t^\top - \underline{\mathbf{Z}}_t \underline{\mathbf{Z}}_t^\top)$ will align to $\hat{\Delta}$ perfectly. Next, we can derive the projection matrix for the top- $2r^*$ subspace of $\frac{1}{2} (\mathbf{Z}_t \mathbf{Z}_t^\top - \underline{\mathbf{Z}}_t \underline{\mathbf{Z}}_t^\top)$. First, we have

$$\mathbf{Z}_t \mathcal{P}^\top(\mathbf{Z}_t) = \begin{bmatrix} \mathbf{A}_t (\mathbf{A}_t^\top \mathbf{A}_t)^{-1} \mathbf{A}_t^\top & \mathbf{A}_t (\mathbf{B}_t \mathbf{B}_t^\top)^{-1} \mathbf{B}_t^\top \\ \mathbf{B}_t^\top (\mathbf{A}_t^\top \mathbf{A}_t)^{-1} \mathbf{A}_t^\top & \mathbf{B}_t^\top (\mathbf{B}_t \mathbf{B}_t^\top)^{-1} \mathbf{B}_t^\top \end{bmatrix},$$

which can imply

$$\begin{aligned} \frac{1}{2} \mathbf{Z}_t \mathcal{P}^\top(\mathbf{Z}_t) \mathbf{Z}_t \mathbf{Z}_t^\top &= \frac{1}{2} \begin{bmatrix} \mathbf{A}_t (\mathbf{A}_t^\top \mathbf{A}_t)^{-1} \mathbf{A}_t^\top & \mathbf{A}_t (\mathbf{B}_t \mathbf{B}_t^\top)^{-1} \mathbf{B}_t^\top \\ \mathbf{B}_t^\top (\mathbf{A}_t^\top \mathbf{A}_t)^{-1} \mathbf{A}_t^\top & \mathbf{B}_t^\top (\mathbf{B}_t \mathbf{B}_t^\top)^{-1} \mathbf{B}_t^\top \end{bmatrix} \begin{bmatrix} \mathbf{A}_t \mathbf{A}_t^\top & \mathbf{A}_t \mathbf{B}_t^\top \\ \mathbf{B}_t^\top \mathbf{A}_t^\top & \mathbf{B}_t^\top \mathbf{B}_t^\top \end{bmatrix} \\ &= \frac{1}{2} \begin{bmatrix} \mathbf{A}_t \mathbf{A}_t^\top & \mathbf{A}_t \mathbf{B}_t^\top \\ \mathbf{B}_t^\top \mathbf{A}_t^\top & \mathbf{B}_t^\top \mathbf{B}_t^\top \end{bmatrix} = \frac{1}{2} \mathbf{Z}_t \mathbf{Z}_t^\top. \end{aligned}$$

Similarly, we can derive

$$\frac{1}{2} \underline{\mathbf{Z}}_t \mathcal{P}^\top(\underline{\mathbf{Z}}_t) \underline{\mathbf{Z}}_t \underline{\mathbf{Z}}_t^\top = \frac{1}{2} \begin{bmatrix} \mathbf{A}_t \mathbf{A}_t^\top & -\mathbf{A}_t \mathbf{B}_t^\top \\ -\mathbf{B}_t^\top \mathbf{A}_t^\top & \mathbf{B}_t^\top \mathbf{B}_t^\top \end{bmatrix} = \frac{1}{2} \underline{\mathbf{Z}}_t \underline{\mathbf{Z}}_t^\top.$$

Additionally, we have

$$\frac{1}{2} \mathbf{Z}_t \mathcal{P}^\top(\underline{\mathbf{Z}}_t) \mathbf{Z}_t \underline{\mathbf{Z}}_t^\top = \mathbf{0}_{(d+k) \times (d+k)}, \quad \frac{1}{2} \underline{\mathbf{Z}}_t \mathcal{P}^\top(\mathbf{Z}_t) \underline{\mathbf{Z}}_t \mathbf{Z}_t^\top = \mathbf{0}_{(d+k) \times (d+k)}.$$

Base on the above identity, we can obtain that the subspace of $\mathbf{Z}_t \mathbf{Z}_t^\top$ is orthogonal to the subspace of $\underline{\mathbf{Z}}_t \underline{\mathbf{Z}}_t^\top$. Since $\text{Rank}(\mathbf{Z}_t \mathbf{Z}_t^\top) \leq r$ and $\text{Rank}(\underline{\mathbf{Z}}_t \underline{\mathbf{Z}}_t^\top) \leq r$, then we have that $\text{Rank}(\mathbf{Z}_t \mathbf{Z}_t^\top - \underline{\mathbf{Z}}_t \underline{\mathbf{Z}}_t^\top) \leq 2r^*$ since $r = r^*$. Therefore, we can construct a valid projection matrix

$$\mathbf{P}_t := \mathbf{Z}_t \mathcal{P}^\top(\mathbf{Z}_t) + \underline{\mathbf{Z}}_t \mathcal{P}^\top(\underline{\mathbf{Z}}_t), \quad (82)$$

which satisfies

$$\frac{1}{2} \mathbf{P}_t (\mathbf{Z}_t \mathbf{Z}_t^\top - \underline{\mathbf{Z}}_t \underline{\mathbf{Z}}_t^\top) = \frac{1}{2} (\mathbf{Z}_t \mathbf{Z}_t^\top - \underline{\mathbf{Z}}_t \underline{\mathbf{Z}}_t^\top),$$

and

$$\frac{1}{2} (\mathbf{I}_{d+k} - \mathbf{P}_t) (\mathbf{Z}_t \mathbf{Z}_t^\top - \underline{\mathbf{Z}}_t \underline{\mathbf{Z}}_t^\top) = \mathbf{0}_{(d+k) \times (d+k)}. \quad (83)$$

Also, we can verify that \mathbf{P}_t is symmetric and $\mathbf{P}_t \mathbf{P}_t = \mathbf{P}_t$. Therefore we can conclude that \mathbf{P}_t is the projection matrix which maps matrices or vectors to the top- $2r$ subspace of $\frac{1}{2} (\mathbf{Z}_t \mathbf{Z}_t^\top - \underline{\mathbf{Z}}_t \underline{\mathbf{Z}}_t^\top)$. For notational simplicity, here we fix the timestamp t and denote

$$\mathbf{F} := \frac{1}{2\sqrt{2}} (\mathbf{Z}_t \mathbf{Z}_t^\top - \underline{\mathbf{Z}}_t \underline{\mathbf{Z}}_t^\top),$$

which means

$$\left\| \mathbf{F} - \frac{\hat{\Delta}}{\sqrt{2}} \right\|_F = \|\mathbf{A}_t \mathbf{B}_t - \Delta\|_F \leq \rho \lambda_{r^*}^*. \quad (84)$$

Next, we define $\mathbf{P}_t := \mathbf{L} \mathbf{L}^\top \in \mathbb{R}^{(d+k) \times (d+k)}$ with

$$\mathbf{L} = \begin{bmatrix} \mathbf{U}_{\mathbf{A}_t} & \mathbf{0}_{d \times r} \\ \mathbf{0}_{k \times r} & \mathbf{V}_{\mathbf{B}_t} \end{bmatrix} \in \mathbb{R}^{(d+k) \times 2r},$$

and $(\mathbf{I}_{d+k} - \mathbf{P}_t) = \mathbf{L}_\perp \mathbf{L}_\perp^\top$ where

$$\mathbf{L}_\perp = \begin{bmatrix} \mathbf{U}_{\mathbf{A}_t, \perp} & \mathbf{0}_{d \times (k-r)} \\ \mathbf{0}_{k \times (d-r)} & \mathbf{V}_{\mathbf{B}_t, \perp} \end{bmatrix} \in \mathbb{R}^{(d+k) \times (d+k-2r)},$$

then we have

$$\begin{aligned} \left\| \mathbf{F} - \frac{\hat{\Delta}}{\sqrt{2}} \right\|_F^2 &= \left\| \begin{bmatrix} \mathbf{L}^\top \\ \mathbf{L}_\perp^\top \end{bmatrix} \left(\mathbf{F} - \frac{\hat{\Delta}}{\sqrt{2}} \right) \begin{bmatrix} \mathbf{L} & \mathbf{L}_\perp \end{bmatrix} \right\|_F^2 \\ &= \left\| \begin{bmatrix} \mathbf{L}^\top \mathbf{F} \mathbf{L} - \mathbf{L}^\top \frac{\hat{\Delta}}{\sqrt{2}} \mathbf{L} & -\mathbf{L}^\top \frac{\hat{\Delta}}{\sqrt{2}} \mathbf{L}_\perp \\ -\mathbf{L}_\perp^\top \frac{\hat{\Delta}}{\sqrt{2}} \mathbf{L} & \mathbf{L}_\perp^\top \frac{\hat{\Delta}}{\sqrt{2}} \mathbf{L}_\perp \end{bmatrix} \right\|_F^2 \quad [\text{by Eq. (83)}] \\ &= \left\| \mathbf{L}^\top \mathbf{F} \mathbf{L} - \frac{1}{\sqrt{2}} \mathbf{L}^\top \hat{\Delta} \mathbf{L} \right\|_F^2 + \frac{1}{2} \left\| \mathbf{L}_\perp^\top \hat{\Delta} \mathbf{L} \right\|_F^2 + \frac{1}{2} \left\| \mathbf{L}^\top \hat{\Delta} \mathbf{L}_\perp \right\|_F^2 + \frac{1}{2} \left\| \mathbf{L}_\perp^\top \hat{\Delta} \mathbf{L}_\perp \right\|_F^2. \quad (85) \end{aligned}$$

Since $\mathbf{I}_{d+k} - \mathbf{P}_t = \mathbf{L}_\perp \mathbf{L}_\perp^\top$, then we have

$$\|(\mathbf{I}_d - \mathbf{U}_{A_t} \mathbf{U}_{A_t}^\top) \Delta (\mathbf{I}_k - \mathbf{V}_{B_t} \mathbf{V}_{B_t}^\top)\|_F = \frac{1}{\sqrt{2}} \|\mathbf{L}_\perp^\top \hat{\Delta} \mathbf{L}_\perp\|_F. \quad (86)$$

Next, by Eq. (85), we have $\|\mathbf{F} - \frac{\hat{\Delta}}{\sqrt{2}}\|_F^2 \geq \frac{1}{2} \|\mathbf{L}_\perp^\top \hat{\Delta} \mathbf{L}_\perp\|_F^2 + \frac{1}{2} \|\mathbf{L}^\top \hat{\Delta} \mathbf{L}_\perp\|_F^2$, leading to

$$\frac{\frac{1}{2} \|\mathbf{L}_\perp^\top \hat{\Delta} \mathbf{L}_\perp\|_F^2}{\|\mathbf{F} - \frac{\hat{\Delta}}{\sqrt{2}}\|_F^2} \leq \frac{\frac{1}{2} \|\mathbf{L}_\perp^\top \hat{\Delta} \mathbf{L}_\perp\|_F^2}{\frac{1}{2} \|\mathbf{L}_\perp^\top \hat{\Delta} \mathbf{L}_\perp\|_F^2 + \frac{1}{2} \|\mathbf{L}^\top \hat{\Delta} \mathbf{L}_\perp\|_F^2}. \quad (87)$$

The technical part is to lower bound $\|\mathbf{L}_\perp^\top \hat{\Delta} \mathbf{L}_\perp\|_F^2$ and $\|\mathbf{L}^\top \hat{\Delta} \mathbf{L}_\perp\|_F^2$, we will rely on the following decomposition which based on Eq. (81), i.e.

$$\begin{aligned} \|\mathbf{L}_\perp^\top \hat{\Delta} \mathbf{L}_\perp\|_F^2 &= \|\mathbf{L}_\perp^\top \hat{\mathbf{U}} \hat{\mathbf{S}} \hat{\mathbf{V}}^\top \mathbf{L}_\perp\|_F^2 \\ &= \left\| \left(\mathbf{L}_\perp^\top \hat{\mathbf{U}} \hat{\mathbf{S}}^{1/2} \right) \left(\mathbf{L}^\top \hat{\mathbf{V}} \hat{\mathbf{S}}^{1/2} \right)^\top \right\|_F^2 \\ &= \text{tr} \left(\left(\mathbf{L}^\top \hat{\mathbf{V}} \hat{\mathbf{S}}^{1/2} \right) \left(\mathbf{L}_\perp^\top \hat{\mathbf{U}} \hat{\mathbf{S}}^{1/2} \right)^\top \left(\mathbf{L}_\perp^\top \hat{\mathbf{U}} \hat{\mathbf{S}}^{1/2} \right) \left(\mathbf{L}^\top \hat{\mathbf{V}} \hat{\mathbf{S}}^{1/2} \right)^\top \right) \\ &= \text{tr} \left(\left(\mathbf{L}^\top \hat{\mathbf{V}} \hat{\mathbf{S}}^{1/2} \right)^\top \left(\mathbf{L}^\top \hat{\mathbf{V}} \hat{\mathbf{S}}^{1/2} \right) \left(\mathbf{L}_\perp^\top \hat{\mathbf{U}} \hat{\mathbf{S}}^{1/2} \right)^\top \left(\mathbf{L}_\perp^\top \hat{\mathbf{U}} \hat{\mathbf{S}}^{1/2} \right) \right) \\ &= \text{tr} \left(\left(\hat{\mathbf{S}}^{1/2} \hat{\mathbf{V}}^\top \mathbf{L} \mathbf{L}^\top \hat{\mathbf{V}} \hat{\mathbf{S}}^{1/2} \right) \left(\hat{\mathbf{S}}^{1/2} \hat{\mathbf{U}}^\top \mathbf{L}_\perp \mathbf{L}_\perp^\top \hat{\mathbf{U}} \hat{\mathbf{S}}^{1/2} \right) \right). \end{aligned}$$

Notice that $\hat{\mathbf{S}}^{1/2} \hat{\mathbf{V}}^\top \mathbf{L} \mathbf{L}^\top \hat{\mathbf{V}} \hat{\mathbf{S}}^{1/2}$ and $\hat{\mathbf{S}}^{1/2} \hat{\mathbf{U}}^\top \mathbf{L}_\perp \mathbf{L}_\perp^\top \hat{\mathbf{U}} \hat{\mathbf{S}}^{1/2}$ are two positive semi-definite matrices, then by lower bound of trace of product of positive semi-definite matrices, using Weyl inequality, we have

$$\begin{aligned} \|\mathbf{L}_\perp^\top \hat{\Delta} \mathbf{L}_\perp\|_F^2 &\geq \lambda_{2r^*} \left(\hat{\mathbf{S}}^{1/2} \hat{\mathbf{V}}^\top \mathbf{L} \mathbf{L}^\top \hat{\mathbf{V}} \hat{\mathbf{S}}^{1/2} \right) \left\| \hat{\mathbf{S}}^{1/2} \hat{\mathbf{U}}^\top \mathbf{L}_\perp \mathbf{L}_\perp^\top \hat{\mathbf{U}} \hat{\mathbf{S}}^{1/2} \right\|_F \\ &\geq \lambda_{r^*}^* \times \lambda_{2r^*} \left(\hat{\mathbf{V}}^\top \mathbf{L} \mathbf{L}^\top \hat{\mathbf{V}} \right) \left\| \hat{\mathbf{S}}^{1/2} \hat{\mathbf{U}}^\top \mathbf{L}_\perp \mathbf{L}_\perp^\top \hat{\mathbf{U}} \hat{\mathbf{S}}^{1/2} \right\|_F \\ &= \lambda_{r^*}^* \times \lambda_{2r^*} \left(\hat{\mathbf{V}}^\top \mathbf{L} \mathbf{L}^\top \hat{\mathbf{V}} \right) \left\| \mathbf{L}_\perp \mathbf{L}_\perp^\top \hat{\mathbf{U}} \hat{\mathbf{S}} \hat{\mathbf{U}}^\top \mathbf{L}_\perp \mathbf{L}_\perp^\top \right\|_F, \end{aligned}$$

where the last equality follows from

$$\begin{aligned} \left\| \hat{\mathbf{S}}^{1/2} \hat{\mathbf{U}}^\top \mathbf{L}_\perp \mathbf{L}_\perp^\top \hat{\mathbf{U}} \hat{\mathbf{S}}^{1/2} \right\|_F^2 &= \text{tr} \left(\hat{\mathbf{S}}^{1/2} \hat{\mathbf{U}}^\top \mathbf{L}_\perp \mathbf{L}_\perp^\top \hat{\mathbf{U}} \hat{\mathbf{S}} \hat{\mathbf{U}}^\top \mathbf{L}_\perp \mathbf{L}_\perp^\top \hat{\mathbf{U}} \hat{\mathbf{S}}^{1/2} \right) \\ &= \text{tr} \left(\hat{\mathbf{S}}^{1/2} \hat{\mathbf{U}}^\top \mathbf{L}_\perp \mathbf{L}_\perp^\top \left(\mathbf{L}_\perp \mathbf{L}_\perp^\top \hat{\mathbf{U}} \hat{\mathbf{S}} \hat{\mathbf{U}}^\top \mathbf{L}_\perp \mathbf{L}_\perp^\top \right) \mathbf{L}_\perp \mathbf{L}_\perp^\top \hat{\mathbf{U}} \hat{\mathbf{S}}^{1/2} \right) \\ &= \text{tr} \left(\left(\mathbf{L}_\perp \mathbf{L}_\perp^\top \hat{\mathbf{U}} \hat{\mathbf{S}} \hat{\mathbf{U}}^\top \mathbf{L}_\perp \mathbf{L}_\perp^\top \right) \left(\mathbf{L}_\perp \mathbf{L}_\perp^\top \hat{\mathbf{U}} \hat{\mathbf{S}} \hat{\mathbf{U}}^\top \mathbf{L}_\perp \mathbf{L}_\perp^\top \right) \right) \\ &= \left\| \mathbf{L}_\perp \mathbf{L}_\perp^\top \hat{\mathbf{U}} \hat{\mathbf{S}} \hat{\mathbf{U}}^\top \mathbf{L}_\perp \mathbf{L}_\perp^\top \right\|_F^2. \end{aligned}$$

Similarly, we have

$$\|\mathbf{L}^\top \hat{\Delta} \mathbf{L}_\perp\|_F^2 \geq \lambda_{r^*}^* \times \lambda_{2r^*} \left(\hat{\mathbf{U}}^\top \mathbf{L} \mathbf{L}^\top \hat{\mathbf{U}} \right) \left\| \mathbf{L}_\perp \mathbf{L}_\perp^\top \hat{\mathbf{V}} \hat{\mathbf{S}} \hat{\mathbf{V}}^\top \mathbf{L}_\perp \mathbf{L}_\perp^\top \right\|_F.$$

Next, we can derive

$$\begin{aligned} \left\| \mathbf{L}_\perp \mathbf{L}_\perp^\top \hat{\mathbf{U}} \hat{\mathbf{S}} \hat{\mathbf{U}}^\top \mathbf{L}_\perp \mathbf{L}_\perp^\top \right\|_F^2 &= \left\| \mathbf{L}_\perp \mathbf{L}_\perp^\top \left(\Phi \mathbf{S}^* \Phi^\top + \underline{\Phi} \mathbf{S}^* \underline{\Phi}^\top \right) \mathbf{L}_\perp \mathbf{L}_\perp^\top \right\|_F^2 && \text{[by Eq. (80)]} \\ &= \left\| \mathbf{L}_\perp \mathbf{L}_\perp^\top \underline{\Phi} \mathbf{S}^* \underline{\Phi}^\top \mathbf{L}_\perp \mathbf{L}_\perp^\top \right\|_F^2 + \left\| \mathbf{L}_\perp \mathbf{L}_\perp^\top \Phi \mathbf{S}^* \Phi^\top \mathbf{L}_\perp \mathbf{L}_\perp^\top \right\|_F^2 \\ &\quad + 2 \left\langle \mathbf{L}_\perp \mathbf{L}_\perp^\top \underline{\Phi} \mathbf{S}^* \underline{\Phi}^\top \mathbf{L}_\perp \mathbf{L}_\perp^\top, \mathbf{L}_\perp \mathbf{L}_\perp^\top \Phi \mathbf{S}^* \Phi^\top \mathbf{L}_\perp \mathbf{L}_\perp^\top \right\rangle, \end{aligned}$$

and

$$\begin{aligned}
 \left\| L_{\perp} L_{\perp}^{\top} \widehat{\mathbf{V}} \widehat{\mathbf{S}} \widehat{\mathbf{V}}^{\top} L_{\perp} L_{\perp}^{\top} \right\|_{\text{F}}^2 &= \left\| L_{\perp} L_{\perp}^{\top} \left(\Phi \mathbf{S}^* \Phi^{\top} + \underline{\Phi} \mathbf{S}^* \underline{\Phi}^{\top} \right) L_{\perp} L_{\perp}^{\top} \right\|_{\text{F}}^2 && \text{[by Eq. (80)]} \\
 &= \left\| L_{\perp} L_{\perp}^{\top} \underline{\Phi} \mathbf{S}^* \underline{\Phi}^{\top} L_{\perp} L_{\perp}^{\top} \right\|_{\text{F}}^2 + \left\| L_{\perp} L_{\perp}^{\top} \Phi \mathbf{S}^* \Phi^{\top} L_{\perp} L_{\perp}^{\top} \right\|_{\text{F}}^2 \\
 &\quad + 2 \left\langle L_{\perp} L_{\perp}^{\top} \underline{\Phi} \mathbf{S}^* \underline{\Phi}^{\top} L_{\perp} L_{\perp}^{\top}, L_{\perp} L_{\perp}^{\top} \Phi \mathbf{S}^* \Phi^{\top} L_{\perp} L_{\perp}^{\top} \right\rangle.
 \end{aligned}$$

Also, we can obtain

$$\begin{aligned}
 \left\| L_{\perp}^{\top} \hat{\Delta} L_{\perp} \right\|_{\text{F}}^2 &= \left\| L_{\perp} L_{\perp}^{\top} \widehat{\mathbf{U}} \widehat{\mathbf{S}} \widehat{\mathbf{U}}^{\top} L_{\perp} L_{\perp}^{\top} \right\|_{\text{F}}^2 \\
 &= \left\| L_{\perp} L_{\perp}^{\top} \left(\Phi \mathbf{S}^* \Phi^{\top} - \underline{\Phi} \mathbf{S}^* \underline{\Phi}^{\top} \right) L_{\perp} L_{\perp}^{\top} \right\|_{\text{F}}^2 && \text{[by Eq. (80)]} \\
 &= \left\| L_{\perp} L_{\perp}^{\top} \underline{\Phi} \mathbf{S}^* \underline{\Phi}^{\top} L_{\perp} L_{\perp}^{\top} \right\|_{\text{F}}^2 + \left\| L_{\perp} L_{\perp}^{\top} \Phi \mathbf{S}^* \Phi^{\top} L_{\perp} L_{\perp}^{\top} \right\|_{\text{F}}^2 - 2 \left\langle L_{\perp} L_{\perp}^{\top} \underline{\Phi} \mathbf{S}^* \underline{\Phi}^{\top} L_{\perp} L_{\perp}^{\top}, L_{\perp} L_{\perp}^{\top} \Phi \mathbf{S}^* \Phi^{\top} L_{\perp} L_{\perp}^{\top} \right\rangle.
 \end{aligned}$$

Notice that the matrix inner product term is the inner product of two positive semi-definite matrices, then by trace inequality for positive semi-definite matrices, we can obtain

$$\left\langle L_{\perp} L_{\perp}^{\top} \underline{\Phi} \mathbf{S}^* \underline{\Phi}^{\top} L_{\perp} L_{\perp}^{\top}, L_{\perp} L_{\perp}^{\top} \Phi \mathbf{S}^* \Phi^{\top} L_{\perp} L_{\perp}^{\top} \right\rangle = \text{tr} \left(\left(L_{\perp} L_{\perp}^{\top} \underline{\Phi} \mathbf{S}^* \underline{\Phi}^{\top} L_{\perp} L_{\perp}^{\top} \right) \left(L_{\perp} L_{\perp}^{\top} \Phi \mathbf{S}^* \Phi^{\top} L_{\perp} L_{\perp}^{\top} \right) \right) \geq 0.$$

Then, we can claim

$$\left\| L_{\perp} L_{\perp}^{\top} \widehat{\mathbf{U}} \widehat{\mathbf{S}} \widehat{\mathbf{U}}^{\top} L_{\perp} L_{\perp}^{\top} \right\|_{\text{F}}^2, \left\| L_{\perp} L_{\perp}^{\top} \widehat{\mathbf{V}} \widehat{\mathbf{S}} \widehat{\mathbf{V}}^{\top} L_{\perp} L_{\perp}^{\top} \right\|_{\text{F}}^2 \geq \left\| L_{\perp}^{\top} \hat{\Delta} L_{\perp} \right\|_{\text{F}}^2. \quad (88)$$

Next, we can obtain

$$\begin{aligned}
 \left\| L_{\perp}^{\top} \hat{\Delta} L_{\perp} \right\|_{\text{F}}^2 + \left\| L_{\perp}^{\top} \hat{\Delta} L_{\perp} \right\|_{\text{F}}^2 &\geq \lambda_{r^*}^* \times \lambda_{2r^*} \left(\widehat{\mathbf{U}}^{\top} L L^{\top} \widehat{\mathbf{U}} \right) \left\| L_{\perp} L_{\perp}^{\top} \widehat{\mathbf{V}} \widehat{\mathbf{S}} \widehat{\mathbf{V}}^{\top} L_{\perp} L_{\perp}^{\top} \right\|_{\text{F}} \\
 &\quad + \lambda_{r^*}^* \times \lambda_{2r^*} \left(\widehat{\mathbf{V}}^{\top} L L^{\top} \widehat{\mathbf{V}} \right) \left\| L_{\perp} L_{\perp}^{\top} \widehat{\mathbf{U}} \widehat{\mathbf{S}} \widehat{\mathbf{U}}^{\top} L_{\perp} L_{\perp}^{\top} \right\|_{\text{F}} \\
 &\geq \lambda_{r^*}^* \min \left\{ \lambda_{2r^*} \left(\widehat{\mathbf{U}}^{\top} L L^{\top} \widehat{\mathbf{U}} \right), \lambda_{2r^*} \left(\widehat{\mathbf{V}}^{\top} L L^{\top} \widehat{\mathbf{V}} \right) \right\} \\
 &\quad \times \left(\left\| L_{\perp} L_{\perp}^{\top} \widehat{\mathbf{V}} \widehat{\mathbf{S}} \widehat{\mathbf{V}}^{\top} L_{\perp} L_{\perp}^{\top} \right\|_{\text{F}} + \left\| L_{\perp} L_{\perp}^{\top} \widehat{\mathbf{U}} \widehat{\mathbf{S}} \widehat{\mathbf{U}}^{\top} L_{\perp} L_{\perp}^{\top} \right\|_{\text{F}} \right) \\
 &\geq 2 \lambda_{r^*}^* \min \left\{ \lambda_{2r^*} \left(\widehat{\mathbf{U}}^{\top} L L^{\top} \widehat{\mathbf{U}} \right), \lambda_{2r^*} \left(\widehat{\mathbf{V}}^{\top} L L^{\top} \widehat{\mathbf{V}} \right) \right\} \left\| L_{\perp}^{\top} \hat{\Delta} L_{\perp} \right\|_{\text{F}}. && \text{[by Eq. (88)]}
 \end{aligned}$$

Then, combining the above inequality and Eq. (87), we have

$$\begin{aligned}
 \frac{\frac{1}{2} \left\| L_{\perp}^{\top} \hat{\Delta} L_{\perp} \right\|_{\text{F}}^2}{\left\| \mathbf{F} - \frac{\hat{\Delta}}{\sqrt{2}} \right\|_{\text{F}}^2} &\leq \frac{\left\| L_{\perp}^{\top} \hat{\Delta} L_{\perp} \right\|_{\text{F}}^2}{2 \lambda_{r^*}^* \min \left\{ \lambda_{2r^*} \left(\widehat{\mathbf{U}}^{\top} L L^{\top} \widehat{\mathbf{U}} \right), \lambda_{2r^*} \left(\widehat{\mathbf{V}}^{\top} L L^{\top} \widehat{\mathbf{V}} \right) \right\} \left\| L_{\perp}^{\top} \hat{\Delta} L_{\perp} \right\|_{\text{F}}} \\
 &= \frac{\left\| L_{\perp}^{\top} \hat{\Delta} L_{\perp} \right\|_{\text{F}}}{2 \lambda_{r^*}^* \min \left\{ \lambda_{2r^*} \left(\widehat{\mathbf{U}}^{\top} L L^{\top} \widehat{\mathbf{U}} \right), \lambda_{2r^*} \left(\widehat{\mathbf{V}}^{\top} L L^{\top} \widehat{\mathbf{V}} \right) \right\}}.
 \end{aligned}$$

Next, we will focus on the lower bound of $\lambda_{2r^*} \left(\widehat{\mathbf{U}}^{\top} L L^{\top} \widehat{\mathbf{U}} \right)$ and $\lambda_{2r^*} \left(\widehat{\mathbf{V}}^{\top} L L^{\top} \widehat{\mathbf{V}} \right)$. Due to symmetry, the technique is identical to each other, so here we only prove for $\lambda_{2r^*} \left(\widehat{\mathbf{U}}^{\top} L L^{\top} \widehat{\mathbf{U}} \right)$. First, $\lambda_{2r^*} \left(\widehat{\mathbf{U}}^{\top} L L^{\top} \widehat{\mathbf{U}} \right) = \lambda_{2r^*}^2 \left(L^{\top} \widehat{\mathbf{U}} \right)$ since $\widehat{\mathbf{U}}^{\top} L L^{\top} \widehat{\mathbf{U}}$ is symmetric. Next, we have

$$\lambda_{2r^*}^2 \left(L^{\top} \widehat{\mathbf{U}} \right) = 1 - \left\| L_{\perp}^{\top} \widehat{\mathbf{U}} \right\|_{\text{op}}^2,$$

where $\left\| \mathbf{L}_\perp^\top \hat{\mathbf{U}} \right\|_{op}$ can be upper bounded by Wedin's $\sin(\Theta)$ theorem, here we use a variant from in Chen et al. (2021b, Theorem 2.9) to obtain

$$\left\| \mathbf{L}_\perp^\top \hat{\mathbf{U}} \right\|_{op} \leq \frac{2 \left\| \mathbf{F} - \frac{\hat{\Delta}}{\sqrt{2}} \right\|_{op}}{\lambda_{2r^*}^* \left(\frac{\hat{\Delta}}{\sqrt{2}} \right)} \leq \frac{2\sqrt{2} \left\| \mathbf{F} - \frac{\hat{\Delta}}{\sqrt{2}} \right\|_F}{\lambda_{r^*}^*} \leq 2\sqrt{2}\rho, \quad \left[\text{by } \left\| \mathbf{F} - \frac{\hat{\Delta}}{\sqrt{2}} \right\|_F \leq \rho\lambda_{r^*}^* \right]$$

which implies

$$\lambda_{2r^*}^2 \left(\mathbf{L}^\top \hat{\mathbf{U}} \right) \geq 1 - 8\rho^2. \quad (89)$$

Therefore, we have

$$\begin{aligned} \rho^2(\lambda_{r^*}^*)^2 &= \rho^2 \lambda_{2r^*}^2(\hat{\Delta}) \geq \left\| \mathbf{F} - \frac{\hat{\Delta}}{\sqrt{2}} \right\|_F^2 \\ &\geq \frac{1}{2} \left\| \mathbf{L}_\perp^\top \hat{\Delta} \mathbf{L} \right\|_F^2 + \frac{1}{2} \left\| \mathbf{L}^\top \hat{\Delta} \mathbf{L}_\perp \right\|_F^2 \quad [\text{by Eq. (85)}] \\ &\geq \lambda_{r^*}^* \min \left\{ \lambda_{2r^*} \left(\hat{\mathbf{U}}^\top \mathbf{L} \mathbf{L}^\top \hat{\mathbf{U}} \right), \lambda_{2r^*} \left(\hat{\mathbf{V}}^\top \mathbf{L} \mathbf{L}^\top \hat{\mathbf{V}} \right) \right\} \left\| \mathbf{L}_\perp^\top \hat{\Delta} \mathbf{L}_\perp \right\|_F \\ &\geq \frac{1}{2} \lambda_{r^*}^* \left\| \mathbf{L}_\perp^\top \hat{\Delta} \mathbf{L}_\perp \right\|_F, \quad [\text{by Eq. (89) and } \rho \leq 1/4] \end{aligned}$$

which implies

$$\frac{\left\| \mathbf{L}_\perp^\top \hat{\Delta} \mathbf{L}_\perp \right\|_F}{\lambda_{r^*}^*} \leq 2\rho^2.$$

Finally, combining Eq. (86), we can obtain

$$\left\| (\mathbf{I}_d - \mathbf{U}_{\mathbf{A}_t} \mathbf{U}_{\mathbf{A}_t}^\top) \Delta (\mathbf{I}_k - \mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^\top) \right\|_F^2 = \frac{1}{2} \left\| \mathbf{L}_\perp^\top \hat{\Delta} \mathbf{L}_\perp \right\|_F^2 \leq \frac{\rho^2}{1 - 8\rho^2} \left\| \mathbf{A}_t \mathbf{B}_t - \Delta \right\|_F^2.$$

Notice that

$$\begin{aligned} \frac{1}{2} \left\| \mathbf{L}_\perp^\top \hat{\Delta} \mathbf{L} \right\|_F^2 + \frac{1}{2} \left\| \mathbf{L}^\top \hat{\Delta} \mathbf{L}_\perp \right\|_F^2 &= \frac{1}{2} \left\| \mathbf{P}_t \hat{\Delta} (\mathbf{I}_{d+k} - \mathbf{P}_t) \right\|_F^2 + \frac{1}{2} \left\| (\mathbf{I}_{d+k} - \mathbf{P}_t) \hat{\Delta} \mathbf{P}_t \right\|_F^2 \\ &= \left\| (\mathbf{I}_d - \mathbf{U}_{\mathbf{A}_t} \mathbf{U}_{\mathbf{A}_t}^\top) \Delta \mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^\top \right\|_F^2 + \left\| \mathbf{U}_{\mathbf{A}_t} \mathbf{U}_{\mathbf{A}_t}^\top \Delta (\mathbf{I}_k - \mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^\top) \right\|_F^2 \\ &= \left\| (\mathbf{I}_d - \mathbf{U}_{\mathbf{A}_t} \mathbf{U}_{\mathbf{A}_t}^\top) \Delta \mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^\top + \mathbf{U}_{\mathbf{A}_t} \mathbf{U}_{\mathbf{A}_t}^\top \Delta (\mathbf{I}_k - \mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^\top) \right\|_F^2, \end{aligned}$$

then by the decomposition in Eq. (85), we can obtain

$$\left\| (\mathbf{I}_d - \mathbf{U}_{\mathbf{A}_t} \mathbf{U}_{\mathbf{A}_t}^\top) \Delta \mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^\top + \mathbf{U}_{\mathbf{A}_t} \mathbf{U}_{\mathbf{A}_t}^\top \Delta (\mathbf{I}_k - \mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^\top) \right\|_F^2 \leq \left\| \mathbf{F} - \frac{\hat{\Delta}}{\sqrt{2}} \right\|_F^2 = \left\| \mathbf{A}_t \mathbf{B}_t - \Delta \right\|_F^2,$$

which completes the proof. \square

Based on the above estimation, we are ready to deliver the linear convergence rate of $\left\| \mathbf{A}_t \mathbf{B}_t - \Delta \right\|_F$.

Theorem D.10. Suppose $\epsilon \leq \frac{\rho}{3C^* K^2 \gamma \sqrt{2dr^* \kappa}}$ for a positive constant $\rho \leq \frac{1}{20}$, we take $\gamma = 2$ for (Spectral-init), set $\eta \in (c_\eta, 1)$ where $c_\eta > 0$ is a small constant, under assumptions in Section 2.1 for the nonlinear setting and Assumption 4.1, then with probability at least $1 - 2Cdk \exp(-\epsilon^2 N)$ for a universal constant $C > 0$, we have

$$\left\| \mathbf{A}_t \mathbf{B}_t - \Delta \right\|_F \leq \left(1 - \frac{\eta}{4} \right)^t \rho \lambda_{r^*}^*.$$

Proof. We prove it by induction. The following hypothesis holds at $t = 0$ by Lemma D.5, i.e.

$$\|\mathbf{A}_t \mathbf{B}_t - \Delta\|_F \leq \rho \lambda_{r^*}^*.$$

We suppose it also holds for at time t , then the conditions of Lemma D.7, Lemma D.8, and Lemma D.9 are fulfilled. By consequence, we can show that

$$\lambda_{r^*}(\mathbf{A}_t \mathbf{B}_t) \geq (1 - \rho) \lambda_{r^*}^*. \quad [\text{by Weyl's inequality}]$$

Next, by Eq. (77) from Lemma D.6, under initial conditions from Lemma D.5, for time $t + 1$, we can derive

$$\begin{aligned} & \|\mathbf{A}_{t+1} \mathbf{B}_{t+1} - \Delta\|_F \\ & \leq (1 - \eta) \|\mathbf{U}_{\mathbf{A}_t} \mathbf{U}_{\mathbf{A}_t}^\top (\mathbf{A}_t \mathbf{B}_t - \Delta) \mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^\top\|_F \\ & \quad + (1 - \eta/2) \|(\mathbf{I}_d - \mathbf{U}_{\mathbf{A}_t} \mathbf{U}_{\mathbf{A}_t}^\top) (\mathbf{A}_t \mathbf{B}_t - \Delta) \mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^\top + \mathbf{U}_{\mathbf{A}_t} \mathbf{U}_{\mathbf{A}_t}^\top (\mathbf{A}_t \mathbf{B}_t - \Delta) (\mathbf{I}_k - \mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^\top)\|_F \\ & \quad + \|(\mathbf{I}_d - \mathbf{U}_{\mathbf{A}_t} \mathbf{U}_{\mathbf{A}_t}^\top) (\mathbf{A}_t \mathbf{B}_t - \Delta) (\mathbf{I}_k - \mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^\top)\|_F \\ & \quad + 2\eta \|\Xi_t\|_F + \eta^2 \|\mathbf{J}_{\mathbf{W}_t} \mathcal{V}_t \mathcal{S}_t^{-1} \mathcal{U}_t^\top \mathbf{J}_{\mathbf{W}_t}\|_F \\ & \leq (1 - \eta) \|\mathbf{A}_t \mathbf{B}_t - \Delta\|_F \\ & \quad + \left(1 - \eta/2 + \frac{\rho}{\sqrt{1 - 8\rho^2}}\right) \|\mathbf{A}_t \mathbf{B}_t - \Delta\|_F \quad [\text{by Lemma D.9}] \\ & \quad + 2\eta \left(\mathcal{O}\left(\frac{1}{\kappa r^*}\right) + C^* K^2 \sqrt{d} \epsilon\right) \|\mathbf{A}_t \mathbf{B}_t - \Delta\|_F \quad [\text{by Lemma D.7}] \\ & \quad + \eta^2 (1 + \epsilon)^2 \frac{\rho}{1 - \rho} \|\mathbf{A}_t \mathbf{B}_t - \Delta\|_F \quad [\text{by Lemma D.8}] \\ & \leq \left(2 - 3\eta/2 + \frac{\rho}{\sqrt{1 - 8\rho^2}}\right) \|\mathbf{A}_t \mathbf{B}_t - \Delta\|_F \\ & \quad + \eta \left(\frac{2\rho}{3\gamma\sqrt{2}r^*\kappa}\right) \|\mathbf{A}_t \mathbf{B}_t - \Delta\|_F \\ & \quad + \eta^2 \left(1 + \frac{\rho}{3C^* K^2 \gamma \sqrt{2d} r^* \kappa}\right)^2 \frac{\rho}{1 - \rho} \|\mathbf{A}_t \mathbf{B}_t - \Delta\|_F, \quad \left[\text{since } \epsilon \leq \frac{\rho}{3C^* K^2 \gamma \sqrt{2d} r^* \kappa}\right] \end{aligned}$$

with probability at least $1 - 2Cdk \exp(-\epsilon^2 N)$ for a universal constant $C > 0$. Since we take $\rho \leq \frac{1}{20}$, and $\frac{\rho}{\sqrt{1 - 8\rho^2}}$ is monotonically increasing, then there exists a constant $c_\eta > 0$ such that $\forall \eta \in (c_\eta, 1)$, we have

$$\|\mathbf{A}_{t+1} \mathbf{B}_{t+1} - \Delta\|_F \leq \left(1 - \frac{\eta}{4}\right) \|\mathbf{A}_t \mathbf{B}_t - \Delta\|_F.$$

Then, we can obtain the inductive hypothesis at $t + 1$ and prove the claim. \square

E. Auxiliary Results for Proofs

In this subsection, we present some auxiliary results that are needed for our proof. First, we present the estimation of the spectral norm of random matrices. It can be easily derived from (Vershynin, 2018) and we put it here for the completeness.

Lemma E.1. (Vershynin, 2018, Adapted from Theorem 4.6.1) For a random sub-Gaussian matrix $\widetilde{\mathbf{X}} \in \mathbb{R}^{N \times d}$ whose rows are i.i.d. isotropic sub-gaussian random vector with sub-Gaussian norm K , then we have the following statement

$$\mathbb{P}\left(\left\|\frac{1}{N} \widetilde{\mathbf{X}}^\top \widetilde{\mathbf{X}} - \mathbf{I}_d\right\|_{op} > \delta\right) \leq 2 \exp(-CN \min(\delta^2, \delta)).$$

for a universal constant C depending only on K .

Lemma E.2. (*Vershynin, 2010, Adapted from Corollary 5.35*) For a random standard Gaussian matrix $\mathbf{S} \in \mathbb{R}^{d \times r}$ with $[\mathbf{S}]_{ij} \sim \mathcal{N}(0, 1)$, if $d > 2r$, we have

$$\frac{\sqrt{d}}{2} \leq \|\mathbf{S}\|_{op} \leq (2\sqrt{d} + \sqrt{r}), \quad (90)$$

with probability at least $1 - C \exp(-d)$ for some positive constants C .

The following results are modified from the proof of [Stöger & Soltanolkotabi \(2021, Lemma 8.7\)](#).

Lemma E.3. Suppose $\mathbf{S} \in \mathbb{R}^{d \times r}$ is a random standard Gaussian matrix with $[\mathbf{S}]_{ij} \sim \mathcal{N}(0, 1)$ and $\mathbf{U} \in \mathbb{R}^{d \times r^*}$ has orthonormal columns. If $r \geq 2r^*$, with probability at least $1 - C \exp(-r)$ for some positive constants C , we have

$$\lambda_{\min}(\mathbf{U}^\top \mathbf{S}) \gtrsim 1.$$

If $r^* \leq r < 2r^*$, by choosing $\xi > 0$ appropriately, with probability at least $1 - (C\xi)^{r-r^*+1} - C' \exp(-r)$ for some positive constants C, C' , we have

$$\lambda_{\min}(\mathbf{U}^\top \mathbf{S}) \gtrsim \frac{\xi}{r}.$$

Lemma E.4. (*Brutzkus & Globerson, 2017, Lemma 3.2*) Define $\theta(\mathbf{w}, \mathbf{v}) = \cos^{-1} \left(\frac{\langle \mathbf{w}, \mathbf{v} \rangle}{\|\mathbf{w}\|_2 \|\mathbf{v}\|_2} \right)$, then we have

$$h(\mathbf{w}, \mathbf{v}) := \frac{\partial}{\partial \mathbf{w}} \mathbb{E}_{\tilde{\mathbf{x}} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)} [\sigma(\langle \mathbf{w}, \tilde{\mathbf{x}} \rangle) \sigma(\langle \mathbf{v}, \tilde{\mathbf{x}} \rangle)] = \frac{1}{2\pi} \left[\frac{\|\mathbf{v}\|_2}{\|\mathbf{w}\|_2} \sin \theta(\mathbf{w}, \mathbf{v}) \mathbf{w} + (\pi - \theta(\mathbf{w}, \mathbf{v})) \mathbf{v} \right].$$

F. Detailed Comparison with LoRA-GA

LoRA-GA proposes the following initialization strategy

$$\begin{aligned} \mathbf{A}_0 &= -\frac{k^{1/4}}{c} [\tilde{\mathbf{U}}_{\mathbf{G}^\natural}]_{[:,1:r]}, \mathbf{B}_0 = \frac{k^{1/4}}{c} [\tilde{\mathbf{V}}_{\mathbf{G}^\natural}]_{[:,r+1:2r]}^\top, \\ \mathbf{W}_{\text{off}}^\natural &:= \mathbf{W}^\natural - \frac{\alpha}{\sqrt{r}} \mathbf{A}_0 \mathbf{B}_0, \end{aligned} \quad (91)$$

where k is the output dimension, c is a user-specified hyperparameter in constant order. They propose to recover the one-step full gradient \mathbf{G}^\natural to the largest extent after the first LoRA update, i.e., under gradient descent with stepsize η , the adapted weight becomes

$$\begin{aligned} \mathbf{W}_{\text{off}}^\natural + \mathbf{A}_1 \mathbf{B}_1 &:= \mathbf{W}_{\text{off}}^\natural + \frac{\alpha}{\sqrt{r}} \mathbf{A}_0 \mathbf{B}_0 + \frac{\alpha}{\sqrt{r}} \left[-\eta \mathbf{G}^\natural \mathbf{B}_0^\top \mathbf{B}_0 - \eta \mathbf{A}_0 \mathbf{A}_0^\top \mathbf{G}^\natural + \eta^2 \mathbf{G}^\natural \mathbf{B}_0^\top \mathbf{A}_0^\top \mathbf{G}^\natural \right] \\ &= \mathbf{W}^\natural + \frac{\alpha}{\sqrt{r}} \underbrace{\left[-\eta \mathbf{G}^\natural \mathbf{B}_0^\top \mathbf{B}_0 - \eta \mathbf{A}_0 \mathbf{A}_0^\top \mathbf{G}^\natural + \eta^2 \mathbf{G}^\natural \mathbf{B}_0^\top \mathbf{A}_0^\top \mathbf{G}^\natural \right]}_{\text{update in the full parameter space}}. \end{aligned}$$

Then, \mathbf{A}_0 and \mathbf{B}_0 in Eq. (91) can admit the best rank- $2r$ approximation of \mathbf{G}^\natural in terms of full parameter update as they drop the η^2 -term. However, this scheme has structural limitations in various perspectives.

First, as pointed by our theory, \mathbf{B}_t will align to the right-side rank- r^* singular subspace of \mathbf{G}^\natural under random initialization. That means, due to the way LoRA-GA chooses the $(r+1)$ -th to $2r$ -th singular values for \mathbf{B}_0 , the iterate \mathbf{B}_t does not lie in the desired subspace and may not escape an undesirable subspace.

Second, for common LoRA-based algorithms including ours, the global optimization problem is to solve

$$\min_{\mathbf{A}, \mathbf{B}} \left\| \frac{\alpha}{\sqrt{r}} \mathbf{A} \mathbf{B} - \Delta \right\|_{\text{F}}^2, \quad (92)$$

which achieves the global minimum at the best rank- r approximation of Δ . However, under LoRA-GA in Eq. (91), the modifications to the pre-trained weight, i.e. \mathbf{W}^\natural , lead to an unfavorable optimization problem, i.e. to solve

$$\begin{aligned} & \min_{\mathbf{A}, \mathbf{B}} \left\| \mathbf{W}^\natural + \frac{\alpha}{\sqrt{r}} (\mathbf{A}\mathbf{B} - \mathbf{A}_0\mathbf{B}_0) - \widetilde{\mathbf{W}} \right\|_F^2 \\ \Leftrightarrow & \min_{\mathbf{A}, \mathbf{B}} \left\| \frac{\alpha}{\sqrt{r}} \mathbf{A}\mathbf{B} - \left(\frac{\alpha}{\sqrt{r}} \mathbf{A}_0\mathbf{B}_0 + \Delta \right) \right\|_F^2, \end{aligned}$$

which achieves the global minimum at a **biased** best rank- r approximation of Δ , i.e. $\frac{\alpha}{\sqrt{r}} \mathbf{A}_0\mathbf{B}_0 + \Delta$, **no matter what initialization it is**. This upward bias scales in the order of $\Theta(\sqrt{k})$ as they propose the scaling to be \sqrt{k}/c^2 for stability and can be dominant if it has stronger signal than downstream feature Δ .

Lastly, since they ignore η^2 -term in the illustrative analysis, this imposes a latent assumption that the best rank- $2r$ approximation of \mathbf{G}^\natural in terms of full parameter update only holds if the stepsize $\eta \approx 0$. This restriction is consistent with ablation results from Wang et al. (2024) that LoRA-GA is not robust under moderate/large stepsize.

In contrast, LoRA-One aligns well with our theory under correct subspace specification. We do not modify the pre-trained weight so the optimization problem remain the same as Eq. (92). Also, our method is robust to the choice of stepsizes and can undertake large stepsize to achieve faster convergence as shown in the Appendix G.2.

G. Experimental Settings and Additional Results

In Appendix G.1, we firstly provide the experimental details of small-scale experiments in our text. Experimental settings of various NLP tasks in the main text are given by Appendix G.2, Appendix G.3, and Appendix G.4, respectively. Finally, we visualize the spectral properties of both the pre-trained weights and the difference weights after fine-tuning in Appendix G.5 to justify the validity of Assumption 4.1. All small-scale experiments were performed on AMD EPYC 7B12 CPU. All experiments for T5 base model and Llama 2-7B were performed on Nvidia A100 GPUs (40GB).

G.1. Small-Scale Experiments

Here we give the experimental details of Fig. 1(c), Fig. 2, Fig. 3, and Fig. 4, respectively.

Details for Fig. 1(c) The experimental settings are sourced from Meng et al. (2024). We use 10000 odd-labeled data from MNIST (LeCun, 1998) for pre-training and 1000 even-labeled data for fine-tuning. The learning rates for Full Fine-tuning, LoRA, and LoRA-One are set to 5×10^{-4} .

Details for Fig. 2: We initialize \mathbf{A}_0 and \mathbf{B}_0 via (LoRA-init) over variance $\alpha^2 \in \{1, 0.1, 0.01, 0.001, 0.0001\}$. We examine for dimension $d = k = 100$ and $d = k = 1000$. We set $N = 16d$, $r^* = 4$, and $r = 8$. We construct $\Delta := \mathbf{U}\mathbf{V}^\top$ where $\mathbf{U} \in \mathbb{R}^{100 \times 4}$ and $\mathbf{V} \in \mathbb{R}^{100 \times 4}$ are obtained from the SVD of a matrix whose elements are independently sampled from $\mathcal{N}(0, 1)$. We set learning rate $\eta = \frac{1}{64}$. We run 1500 GD steps for each case.

Details for Fig. 4: We take $d = k = 100$ and $N = 12800$ in common. For: 1) under-ranked case $r = 4, r^* = 8$, 2) over-ranked case $r = 8, r^* = 4$. We sample each element of \mathbf{W}^\natural independently from $\mathcal{N}(0, 1)$. We construct $\Delta := \mathbf{U}\mathbf{S}\mathbf{V}^\top$ where $\mathbf{U} \in \mathbb{R}^{100 \times r^*}$ and $\mathbf{V} \in \mathbb{R}^{100 \times r^*}$ are obtained from the SVD of a matrix whose elements are independently sampled from $\mathcal{N}(0, 1)$ and the diagonal values of \mathbf{S} is the first r^* elements of the dictionary $\{40, 30, 20, 10, 1, 1, 1, 0.5\}$. For LoRA-GA defined in Eq. (91), we use learning rate $\eta = 0.5$ and stable parameter 16. For LoRA-SB and LoRA-One, we use learning rate $\eta = 0.5$ and scaling parameter 1.

Comparison on GD trajectories of Fig. 3: Here we conduct a toy experiment to intuitively compare the GD trajectories under (Spectral-init) and (LoRA-init). We fine-tune a simple pre-trained model $y = \mathbf{x}^\top \mathbf{w}^\natural$ on downstream data generated by $\tilde{y} = \tilde{\mathbf{x}}^\top (\mathbf{w}^\natural + \mathbf{w})$, where $\mathbf{x}^\top, \tilde{\mathbf{x}}, \mathbf{w}^\natural, \mathbf{w} \in \mathbb{R}^2$ and $y, \tilde{y} \in \mathbb{R}$. We propose to use LoRA to fine-tune this model by $\hat{y} = \tilde{\mathbf{x}}^\top (\mathbf{w}^\natural + b\mathbf{a})$ where $\mathbf{a} = [a_1, a_2]^\top \in \mathbb{R}^2$ and $b \in \mathbb{R}$. Without loss of generality, we set $\mathbf{w}^\natural = \mathbf{0}$ and $\mathbf{w} = [2, 1]^\top$. The set of global minimizers to this problem is $\{a_1^* = 2/t, a_2^* = 1/t, b^* = t \mid t \in \mathbb{R}\}$. We generate 4 data points $(\tilde{\mathbf{x}}_1, \tilde{\mathbf{x}}_2, \tilde{\mathbf{x}}_3, \tilde{\mathbf{x}}_4)$ whose elements are independently sampled from $\mathcal{N}(0, 1)$ and calculate for $(\tilde{y}_1, \tilde{y}_2, \tilde{y}_3, \tilde{y}_4)$. We use the squared loss $\frac{1}{8} \sum_{i=1}^4 (\tilde{y}_i - b\tilde{\mathbf{x}}_i^\top \mathbf{a})^2$. For (LoRA-init), we initialize each element of \mathbf{a}_0 from $\mathcal{N}(0, 1)$ and $b_0 = 0$. Notice that the variance 1 follows from the Kaiming initialization (He et al., 2015). For (Spectral-init), we first calculate the one-step full gradient, i.e. $\mathbf{g}^\natural := \frac{1}{4} \sum_{i=1}^4 \tilde{y}_i^2 \tilde{\mathbf{x}}_i$. Accordingly, we initialize $\mathbf{a}_0 = \frac{\mathbf{g}^\natural}{\sqrt{\|\mathbf{g}^\natural\|_2}}$ and $b_0 = \sqrt{\|\mathbf{g}^\natural\|_2}$. Next, we run GD to train \mathbf{a} and b for

1000 steps with learning rate $\eta = 0.1$. For each initialization strategy and data generation, we run for 2 different seeds.

G.2. Natural Language Generation

The common hyperparameters are presented in Table 6. Next, we present in the order of {MetaMathQA, Code-Feedback, Alpaca}. We search the best learning rate over $\{5 \times 10^{-4}, 2 \times 10^{-4}, 1 \times 10^{-4}, 5 \times 10^{-5}, 2 \times 10^{-5}, 1 \times 10^{-5}\}$ and batch size over $\{16, 32, 128\}$. The optimized learning rate and batch size are presented in Table 7. Additionally, the scale parameters are set to $\{128, 16, 128\}$ for LoRA-One and $\{64, 64, 64\}$ for LoRA-GA.

Furthermore, we employ the gradient approximation approach proposed by Wang et al. (2024) to replace the full-batch full gradient with stochastic full gradient using a smaller sampled batch from training data and denote the corresponding sample size as Gradient Batch Size. According to the ablation studies on spectral properties and performance under various gradient batch sizes in (Wang et al., 2024), larger gradient batch size only can yield marginal improvement, indicating that it is sufficient to use a smaller batch size for computational efficiency.

Table 6: Common hyperparameters for fine-tuning LLaMA 2-7B on MetaMathQA, Code-Feedback, and Alpaca.

Epoch	Optimizer	(β_1, β_2)	ϵ	LoRA Precision	Weight Decay
1	AdamW	(0.9, 0.999)	1×10^{-8}	FP32	0
Warm-up Ratio	LoRA α	LR Scheduler	Max Length	#Runs	Gradient Batch Size
0.03	16	cosine	1024	3	8

Table 7: Optimized hyperparameters for LoRA, LoRA-GA, and LoRA-One.

	Batch Size	Learning Rate
LoRA	$\{32, 32, 32\}$	$\{2 \times 10^{-4}, 2 \times 10^{-4}, 5 \times 10^{-5}\}$
LoRA-GA	$\{32, 32, 32\}$	$\{5 \times 10^{-5}, 5 \times 10^{-5}, 1 \times 10^{-5}\}$
LoRA-One	$\{32, 32, 16\}$	$\{2 \times 10^{-4}, 5 \times 10^{-4}, 2 \times 10^{-4}\}$

G.3. Math Reasoning on Full Data and Multiple Epochs

We present the detailed values of Fig. 5 in Table 8. The common hyperparameters are same as Table 6. We search the best learning rate over $\{5 \times 10^{-4}, 2 \times 10^{-4}, 1 \times 10^{-4}, 5 \times 10^{-5}, 2 \times 10^{-5}, 1 \times 10^{-5}\}$ and batch size over $\{16, 32, 64, 128\}$. The optimized learning rate and batch size are presented in Table 9. Additionally, the scale parameter are set to 128 for LoRA-One and 64 for LoRA-GA. The imbalance parameter for LoRA+ is set to 16. The results of LoRA, LoRA+, and LoRA-GA are taken from (Wang et al., 2024) since their optimized hyperparameters align with our search.

Table 8: Performance comparison across different methods and epochs

	Epoch 1	Epoch 2	Epoch 3	Epoch 4
LoRA	55.19	58.37	59.28	58.90
LoRA+	56.37	59.21	59.93	59.97
LoRA-GA	56.48	58.64	60.16	60.88
LoRA-One	57.54	60.84	62.62	63.80

G.4. Natural Language Understanding

In Section 6.1, we have presented the experimental comparisons between Algorithm 1 and typical LoRA based algorithms. We follow the configuration of prompt tuning as Wang et al. (2024). The general hyperparameter settings are provides in Table 10. To ensure a fair comparison, we tune the learning rate via grid search over $\{1 \times 10^{-3}, 5 \times 10^{-4}, 2 \times 10^{-4}, 1 \times 10^{-4}\}$.

Table 9: Optimized hyperparameters for LoRA, LoRA+, LoRA-GA, and LoRA-One.

	Batch Size	Learning Rate
LoRA	128	1×10^{-4}
LoRA+	128	5×10^{-5}
LoRA-GA	128	5×10^{-5}
LoRA-One	128	2×10^{-4}

Additionally, the scale parameters for LoRA-One are set to be $\{128, 16, 128, 128, 64\}$ for MNLI, SST-2, CoLA, QNLI, and MRPC.

For Section 6.1, the learning rates for one-step gradient update with gradient batch size 2048 are set to be $\{0.1, 1.0, 0.05, 0.1\}$ for SST-2, CoLA, QNLI, and MRPC. The learning rates for low-rank update ($r = 8$) with gradient batch size 8 are set to be $\{1 \times 10^{-4}, 0.1, 5 \times 10^{-2}, 5 \times 10^{-2}\}$. We omit results for MNLI since the test accuracy remains at 0.0% for the first dozen steps in both full and LoRA fine-tuning, likely due to a substantial structural discrepancy between pre-training and downstream tasks.

Table 10: Common hyperparameters for LoRA fine-tuning on T5-base model.

Epoch	Optimizer	(β_1, β_2)	ϵ	Batch Size	Weight Decay	LR Scheduler
1	AdamW	(0.9, 0.999)	1×10^{-8}	32	0	cosine
Warm-up Ratio	LoRA Alpha	#Runs	Sequence Length	Precision	Gradient Batch Size	
0.03	16	3	128	FP32	8	

G.5. Empirical Verification of Assumption 4.1

We perform full fine-tuning for the pre-trained T5 base model (Raffel et al., 2020) on SST-2 dataset from GLUE (Wang et al., 2019) to approximately access the downstream feature matrices. To ensure better convergence, we take the hyperparameter settings which are presented in Table 11.

To validate the part i) of Assumption 4.1, we inspect the spread of the metric values presented in Fig. 7. We can clearly observe that the ratio between the operator norm of fine-tuned weight and minimum norm of neuron within each layer is bounded.

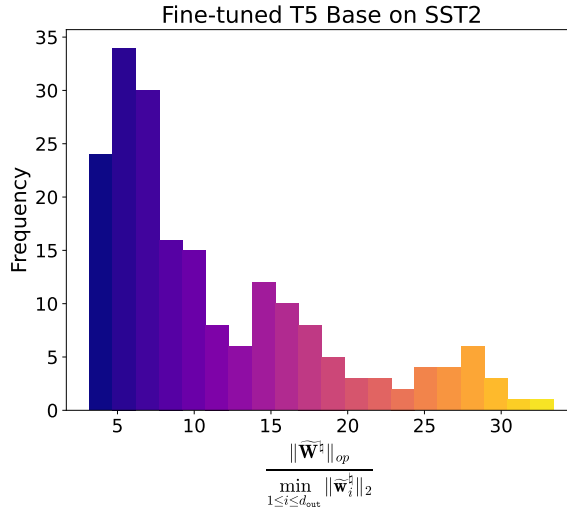


Figure 7: Histogram of the metric values defined presented on x-axis for all fine-tuned weight matrices.

To validate the part ii) of Assumption 4.1, we collect top-32 singular values for each pre-trained layer \mathbf{W}^{\natural} of pre-trained model. After training, we collect top-32 singular values for each difference weights, i.e. $\Delta\mathbf{W} = \mathbf{W}_{\text{fine-tuned}} - \mathbf{W}^{\natural}$. The results are shown in Fig. 8. We observe that, across all layers, the singular values of the pre-trained weights are significantly larger than those of the difference weights. For example, the layer on the right has a pretrained operator norm exceeding 200, while its downstream operator norm is only around 4. Moreover, the singular values decrease drastically as the index increases, indicating an ill-conditioned behavior during fine-tuning.

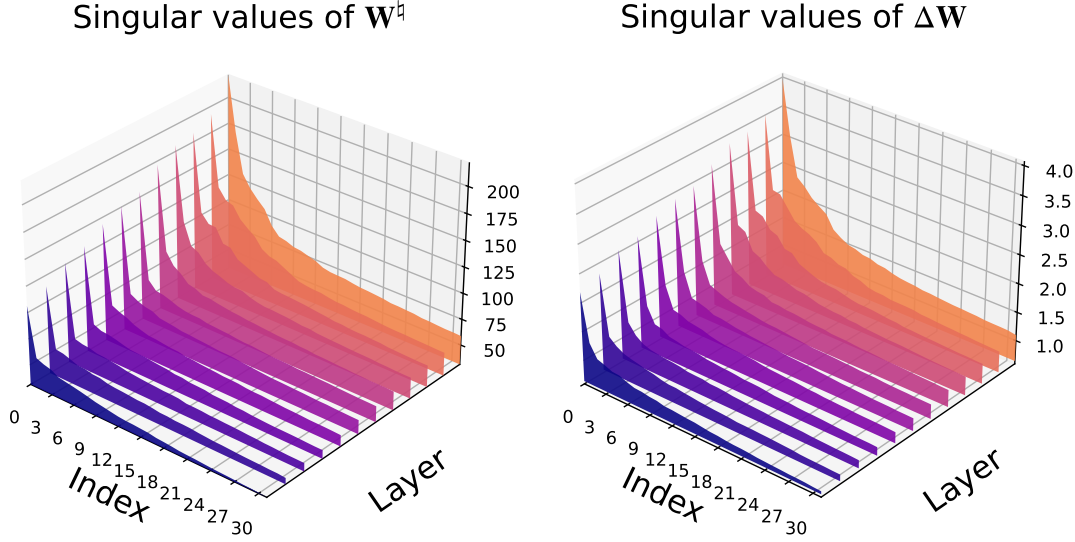


Figure 8: *Left*: top-32 singular values for each pre-trained weight matrices \mathbf{W}^{\natural} . *Right*: top-32 singular values for each difference matrices $\Delta\mathbf{W} = \mathbf{W}_{\text{fine-tuned}} - \mathbf{W}^{\natural}$ after full fine-tuning. The Index is ranked from the largest to the smallest singular values.

Table 11: Hyperparameters for full fine-tuning on T5-base model used for Appendix G.5.

Epoch	Optimizer	(β_1, β_2)	ϵ	Batchsize	Weight Decay	LR	LR Scheduler	Warm-up Ratio
10	AdamW	(0.9, 0.999)	1×10^{-8}	32	0.1	1×10^{-4}	cosine	0.03